



**SOLID STATE
ELECTRONICS
LABORATORY**

JSEP ANNUAL REPORT

1 March, 1995 through 29 February, 1996

**James S. Harris, Jr.
JSEP Principal Investigator
and Program Director**

(415)723-9775

**This work was supported by the
Joint Services Electronics program
(U.S. Army, U.S. Navy and U.S. Air Force)
Contract DAAH04-94-G-0058
and was monitored by the
U.S. Army Research Office**

**Reproduction in whole or in part is permitted
for any purpose of the United States Government**

**This document has been approved for public
release and sale; its distribution is unlimited**

**STANFORD ELECTRONICS LABORATORIES • DEPARTMENT OF ELECTRICAL ENGINEERING
STANFORD UNIVERSITY • STANFORD, CA 94305-4055**

DTIC QUALITY INSPECTED 1

19960909 039

REPORT DOCUMENTATION PAGE			Form Approved OMB No 0704-0188	
<small>Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204 Arlington, VA 22202-4302 and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.</small>				
1. AGENCY USE ONLY (Leave blank)		2. REPORT DATE March 1, 1996	3. REPORT TYPE AND DATES COVERED Annual 1 March 1995 through 29 February 1996	
4. TITLE AND SUBTITLE JSEP Annual Progress Report No. 2			5. FUNDING NUMBERS DAAH04-94-G-0058	
6. AUTHOR(S) J. S. Harris, Program Director				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Stanford University Solid State Electronics Laboratory CIS-X 329 Stanford, CA 94305-4075			8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) U. S. Army Research Office P. O. Box 12211 Research Triangle Park, NC 27709-2211			10. SPONSORING/MONITORING AGENCY REPORT NUMBER ARO 32283.2-EL-JSEP	
11. SUPPLEMENTARY NOTES The view, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy, or decision, unless so designated by other documentation.				
12a. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited.			12b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words)				
14. SUBJECT TERMS			15. NUMBER OF PAGES 69	
			16. PRICE CODE	
17. SECURITY CLASSIFICATION OF REPORT UNCLASSIFIED		18. SECURITY CLASSIFICATION UNCLASSIFIED	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASSIFIED	20. LIMITATION OF ABSTRACT UL

JSEP ANNUAL REPORT

1 March, 1995 through 29 February, 1996

**James S. Harris, Jr.
JSEP Principal Investigator
and Program Director**

(415)723-9775

**This work was supported by the
Joint Services Electronics program
(U.S. Army, U.S. Navy and U.S. Air Force)
Contract DAAH04-94-G-0058
and was monitored by the
U.S. Army Research Office**

**Reproduction in whole or in part is permitted
for any purpose of the United States Government**

**This document has been approved for public
release and sale; its distribution is unlimited**

Abstract

This is the annual report of the research conducted at the Stanford Electronics Laboratories under the sponsorship of the Joint Services Electronics Program from March 1, 1995 through February 29, 1996. This report summarizes the areas of research, identifies the most significant results and lists the dissertations and publications sponsored by contract DAAH04-94-G-0058.

Table of Contents

Introduction and Overview of Principal Accomplishments	3
Unit 1: Investigation of Transport in Quantum Dots	7
Unit 2: Patterned Thin Film Media for High Density Magnetic Recording	15
Unit 3: Investigation of a Metal Source and Drain Field Emission Transistor	20
Unit 4: On-chip Thin Film Solid State Micro-battery	27
Unit 5: CVD Epitaxial Germanium <i>n</i> -channel FETs Formed on Si Substrates using Strain-relief Layers	30
Unit 6: Portable Video on Demand in Wireless Communication	39
Unit 7: Adaptive DFE for GMSK in Indoor Radio Channels	44
Unit 8: Robust Estimation Methods for Adaptive Filtering	57
Unit 9: Efficient Data Compression	64

This work was supported by the Joint Services Electronics Program, contract DAAH04-94-G-0058. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies either expressed or implied of the U.S. Government.

JSEP ANNUAL REPORT

March 1, 1995 - February 29, 1996

Introduction and Overview of Principal Accomplishments

This annual report covers research accomplishments for the period 1 March, 1995 through 29 February, 1996 for basic electronics research conducted in the JSEP program in the Electrical Engineering Department of Stanford University. The Stanford Electronics Lab JSEP Director and Principal Investigator is Professor James Harris. The program work units are as follows:

- Unit 1: Investigation of Transport in Quantum Dots
(James S. Harris)
- Unit 2: Patterned Thin Film Media for High Density Magnetic Recording
(R. Fabian W. Pease)
- Unit 3: Investigation of a Metal Source and Drain Field Emission Transistor
(C. Robert Helms)
- Unit 4: On-chip Thin Film Solid State Micro-battery
(S. Simon Wong)
- Unit 5: CVD Epitaxial Germanium *n*-channel FETs Formed on Si using Strain-relief Layers
(Krishna Saraswat)
- Unit 6: Portable Video on Demand in Wireless Communication
(Teresa H. Y. Meng)
- Unit 7: Adaptive DFE for GMSK in Indoor Radio Channels
(John M. Cioffi)
- Unit 8: Robust Estimation Methods for Adaptive Filtering
(Thomas Kailath)
- Unit 9: Efficient Data Compression
(Thomas M. Cover)

Highlights

In work unit 1, Professor Harris and students have developed the nanofabrication techniques for large (200X200) arrays of 100nm quantum dots and demonstrated the first Coulomb blockade and hysteretic switching behavior in such large arrays. This work represents a significant advance in nanofabrication and demonstrates the robustness of Coulomb blockade compared to quantum interference effects.

In work unit 2, Professor Pease and students have demonstrated and characterized (with Magnetic AFM, alternating gradient magnetometer) magnetic thin film recording media patterned into deep submicron islands for improved density (>12 Gbytes/sq. in.) and lower transition noise. One medium was Polycrystalline Co 20nm thick on Cr which exhibited 1 bit/1 domain/1 island for dimensions less than 150nm. Another medium was single crystal iron film which, when patterned, demonstrated single domain/island behavior for large (1-micron) islands. Magnetic anisotropy in the iron films was dominated by crystalline orientation which allows us to decouple the magnetic direction from the shape of the island; this is valuable for applications involving horizontal recording.

In work unit 5, Professor Saraswat and his students are developing a technology to fabricate high-performance n-channel heterostructure field-effect devices using germanium-rich GeSi grown via graded-alloy strain reduction on (001) silicon substrates. The goal is to combine the high intrinsic electron mobility of germanium and the carrier confinement available with band-structure engineering in the Ge-Si system. Success with the use of graded-alloy epitaxy to isolate the defects associated with the transition between silicon and germanium lattice constants has been achieved.

In work unit 6, Professor Meng's group focuses on low power communications problems. They have developed an energy-on-demand computation system which dynamically adjusts the supply voltage to meet the throughput requirements. A DC-DC switching regulator has been designed that delivers a power-conversion efficiency in excess of 90-percent with tracking speed of under 1 ms. The regulator supplies power ranging from a few mW to several hundreds of mW for all supply voltages of interest. A second effort has produced a real-time, low-power video encoder for pyramid vector quantization (PVQ). This system dissipates only 2.1 mW for real-time video compression of images of 256x256 pixels at 30 frames per second. Applying this quantizer to subband decomposed images, the PVQ encoding delivers better compression performance than the standard JPEG algorithm.

In work unit 9, Professor Cover and students have been investigating the degree to which one can compress images without recognizable distortion. The experiment involves comparison of human vs computer image compression to estimate the minimal rate at which images can be compressed without perceptual distortion. This work is providing a new method for benchmarking data compression algorithms and may lead to a framework to develop entirely new algorithms.

Transfer of Technology

The research results emanating from the JSEP program are usually of either a more fundamental nature or so early that it is not in the vision of more applied programs. Not too surprisingly, such work does not typically lead to instant transfers to industry. However, one

hopes that more fundamental work ultimately has a greater impact because it leads to things that simply would not have been done if left to only research programs with nearer term, clearly identified needs. The transfers of technology described below are thus the result of JSEP supported programs of 5-10 years ago.

Research into the engineering of silicon nanopillars in Professor Pease's JSEP program has led to new insights into the oxidation of silicon under high stress, confined geometry conditions. As Si ULSI continues to shrink, such high stresses are quite important. The results of this research are now being incorporated into SUPREM process models being developed to simulate the processing of next-generation, ultra-small geometry ULSI circuitry.

An essential element in manufacturing high performance AMLCDs is the ability to fabricate TFT driver circuits and integrate them with the liquid crystals on glass substrates. However, the high temperatures and long thermal cycles generally needed to obtain high performance TFTs cause warpage and shrinkage to glass. As a result, fabrication processes are limited to low temperatures and short times. Early work of Professor Krishna Saraswat funded by JSEP and subsequently by DARPA demonstrated high performance TFTs in poly-GeSi with low thermal budget processing, compatible with glass substrates. He demonstrated significantly lower processing temperatures for deposition, doping, recrystallization, and grain boundary passivation. Several novel device structures have been developed to improve TFT performance, such as, increased drive current in the "on" state and reduced leakage in the "off" state. He is actively working with XEROX and Intevac to transfer this technology and several major organizations around the world are now developing the poly-GeSi TFT technology which originated under JSEP support in his laboratory.

The early JSEP work demonstrating the first MBE growth and growth induced layering of the high temperature superconductors by MBE in Professor Harris's program is the basis for the continuing high T_c program at Varian Associates. The focus of their effort is MBE growth induced layering of alternate superconducting and insulating phases to produce well controlled Josephson junctions.

One of the key problems facing modern ultra-high bandwidth communications systems is how to handle the final 100 meters where information delivery is to only a single receiver and the costs of high bandwidth solutions can no longer be divided by a large number of receivers. The early JSEP supported research under Prof. John Cioffi led to the development of the "Discrete MultiTone" (DMT) technology that is now an international standard (ANSI T1.413) for both video transmission and high-speed internet access on twisted pairs, in what is known as Asymmetric Digital Subscriber Lines (ADSL). Stanford holds 4 patents in the area that are exclusively licensed and sublicensed by Stanford to a DMT-spinout, Amati Communications Corporation. Amati has sublicensed the DMT patents to a number of semiconductor and telecom manufacturers around the

world, including Motorola, Northern Telecom, and AT&T (now Lucent Technology). Amati builds products based on the DMT technology and has been extremely successful.

The early JSEP supported work of Professor Tom Cover is now being utilized in many of the WWW browsers. One of the issues is do you wait for all of the information to be supplied serially or do you send information at various levels of refinement so that the description efficiency is optimal at each level ? The idea is to utilize methods of successive refinement to quickly produce a rough picture, then successively more refined pictures until the final version is produced. This work was first described in the paper, "On the successive refinement of information", W. Equitz and T. Cover, IEEE Transactions on Information Theory and is now used by Netscape and many others. Will Equitz was asked by IBM to help write the JPEG data compression standard for progressive transmission based on this work.

UNIT: 1

TITLE: Investigation of Transport in Quantum Dots

PRINCIPAL INVESTIGATOR: J. S. Harris, Jr.

GRADUATE STUDENTS: D. R. Stewart and C. I. Duruöz

1. Scientific Objectives

The continuing drive for increased device density in both IC and memory technologies demands smaller and closer packed future devices. We are pursuing an investigation into the electronic transport in both single quantum devices and large arrays of densely packed quantum dots. A full understanding in both regimes will be required in any successful implementation of single electron electronics. In particular, most studies of quantum devices have concentrated on the very low bias equilibrium behavior [Beenakker][Kouwenhoven]; we concentrate instead on the technologically relevant non-linear high bias operating regime.

We have two main objectives: first, to understand the mechanisms controlling electron transport through single quantum point contacts and quantum dots and second, to study the fundamental characteristics of coulomb blockade and charge coupling in transport through quantum dot arrays.

2. Summary of Research

2.1 Introduction

We previously reported our initial investigations of the electronic transport through 200 x 200 two dimensional quantum dot arrays patterned on a molecular beam epitaxy (MBE) grown GaAs/AlGaAs heterostructure [Harris][Duruöz]. The current-voltage (I-V) relation of the arrays showed two striking features: a threshold for conduction, and multiple switching events accompanied by a hierarchy of hysteresis loops. By changing the voltage applied to a single Schottky gate deposited over the entire array, it was possible to move between the hysteretic and non-hysteretic regimes. A single hysteresis loop was measured in the single control dots fabricated adjacent to the large arrays. No switching or hysteresis was observed above a temperature of 700mK.

We have continued this investigation by focusing on the mechanisms responsible for the switching and hysteresis. It is this behavior, and control of it, that will be most relevant in any technological application.

We have thus characterized in detail the behavior of the single control quantum dots and point contacts in our first generation devices. We have also fabricated a second generation of similar etch defined single devices using a GaAs/AlGaAs heterostructure grown by chemical vapor deposition. All of our single device results have been duplicated in both of these materials to prove the repeatability and robustness of the switching phenomena. Our results show the single hysteresis observed to be the experimental realization of a basic conduction bistability in the I-V relation. When measured on sufficiently fast time scales, the switching bistability manifests as a random telegraph signal in the current under constant voltage bias. Most significantly, we are able to control the bistable switching rate and range with voltages applied to a new back gate and the original front Schottky gate. We are also able to observe the switching in the new devices at a temperature of 4.2K. These results have yielded new insight into the cause of the I-V switching.

2.2 Device Fabrication and Measurement Configuration

All devices measured were fabricated by lithographically patterning a GaAs/AlGaAs epitaxially grown heterostructure. We have utilized a standard modulation doped architecture to create a two dimensional electron gas (2DEG) approximately 800 Å below the wafer surface. First generation and second generation split gate devices were fabricated from MBE material grown in our laboratory with a mobility and sheet density of $\approx 200\,000\text{ cm}^2/\text{Vs}$ and $3.5 \times 10^{11}\text{ cm}^{-2}$; second generation etched devices were patterned on CVD material grown at Sandia National Labs by our collaborator H.Chui with a mobility and density of $\approx 300\,000\text{ cm}^2/\text{Vs}$ and $2.0 \times 10^{11}\text{ cm}^{-2}$.

The devices were formed using electron-beam lithography to define the point contact, dot and array features. Minimum feature size as shown in Fig. 1 is 100 nm, point contact barrier openings are 200-400 nm, and the array periodicity is 800 nm. This lithographic pattern was used as a mask for wet chemical etching 800 Å deep through the 2DEG in the case of etched structures, or NiCrAu metal gate evaporation for the split gate devices. A single 1000 Å Au front gate was deposited over the etched devices. A ground plane below the mounted chips was used as a back gate.

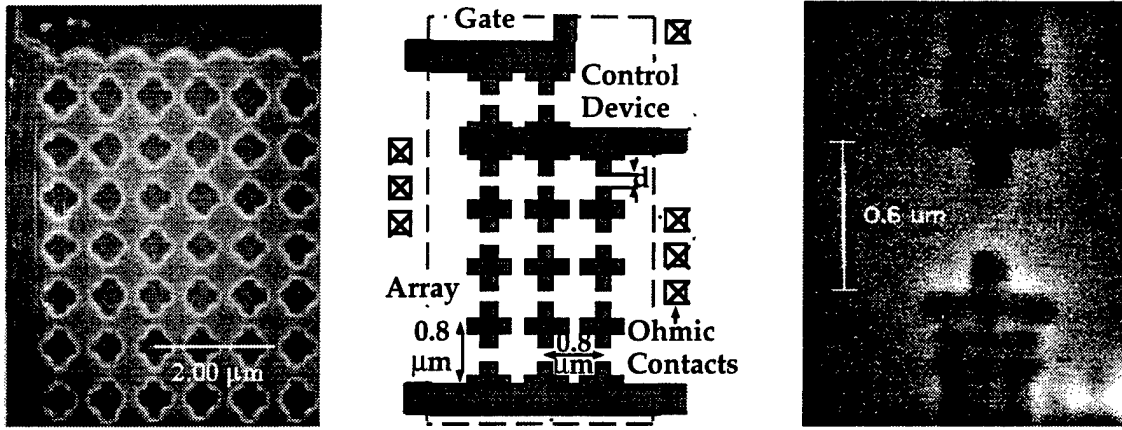


Figure 1: (a) Electron micrograph of part of a 200x200 array. The 2DEG is removed beneath darker regions. (b) Schematic diagram of the array layout (c) Micrograph of a CVD etched point contact of barrier width 260 nm. The 2DEG is removed below the two vertical fingers.

Measurements were conducted at temperatures of 4.2K to 300mK in a pumped He^3 cryostat with a slowly swept dc voltage bias applied to source and drain ohmic contacts across each device. DC voltage biases were also applied to the front and back gates. The dc current was recorded with a slow averaging multimeter and a fast oscilloscope.

2.3 Experimental Results

In Fig. 2 we review the hysteretic behavior observed in the large arrays. All hysteresis loops observed are traversed in a counter clockwise direction in I versus V . This figure shows multiple hysteresis loops shrinking and dissociating into smaller loops as the temperature is raised. All hysteresis has disappeared at 680mK.

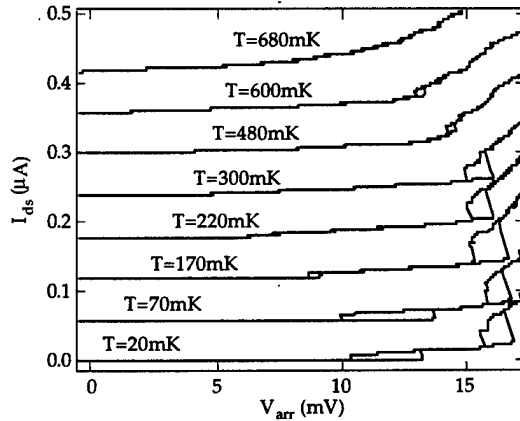


Figure 2: I - V curves (offset for clarity) of an array device at $V_{\text{front gate}} = -115$ mV and various temperatures.

The single point contact and quantum dot devices we have concentrated on all show one bistable switching region as the dc conductance jumps from zero to a finite value, typically $(60 \text{ k}\Omega)^{-1}$. Figure 3 displays how the conductance switches between two bistable states over a small voltage range as the devices turns on. Applying a constant source drain voltage to bias the device at some midpoint of the switching region yields a random telegraph signal in the current as a function of time. (Fig. 3 inset). Recording the times spent in the high and low conductance state yields an average lifetime t_{high} , t_{low} for each state.

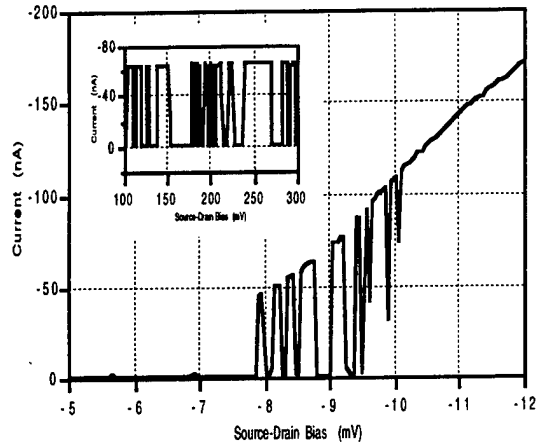


Figure 3: I-V curve of quantum dot displaying bistable conduction switching as the bias is swept up over 8-10 mV. Inset shows random telegraph signal in time at a fixed bias of -9 mV. Temperature is 400 mK.

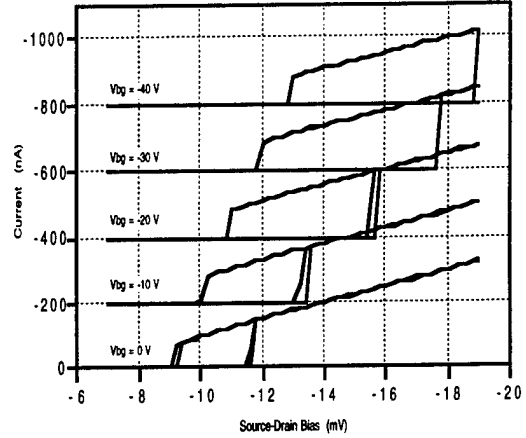


Figure 4: In the hysteric regime, control over the size and position of the hysteresis loop is effected with a backgate voltage as labeled (curves offset for clarity). Results for an etch defined quantum dot at 400 mK.

As the source-drain bias is swept over the switching range these lifetimes appear to change exponentially; t_{high} increases with bias and t_{low} decreases. The clean hysteresis loops initially observed in the arrays can thus be described as bistable conductance regions with average (t_{high} , t_{low}) \gg measurement sweep rate. As the device remains cold, the time constants of this switching increase over several hours until even a slow voltage sweep appears hysteretic.

In this long switching time or 'hysteretic' regime when t_{switch} is much greater than our measurement speed of $O(10\text{s})$, we can use the front and back gates to control the size and position

of the hysteresis. As an increasingly negative backgate voltage is applied, the hysteresis loop expands in size and the initial turn on threshold shifts to higher source-drain bias, as illustrated in Fig. 4.

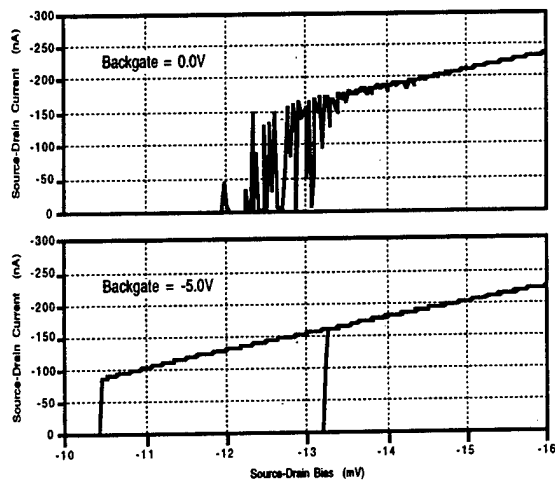


Figure 5: In the fast telegraph switching regime, the backgate is able to reversibly control the bistable state lifetimes. Results from etch defined quantum dot at 400 mK.

In the short switching time or telegraph noise regime we achieve our most significant result; application of a small backgate voltage changes the average state lifetimes dramatically. We are able to continuously control the lifetimes over our full measurement range of $100\mu\text{s}$ to 1000s , seven orders of magnitude. Fig. 5 demonstrates this control.

The CVD etched devices extended the temperature range of this behavior to above 4.2 K. In addition, some of these devices displayed multi-stable switching instead of a simple bistability. The multi-stable devices also showed switching between finite conduction states, and a smoother current turn on. This comparison is made in Fig. 6.

We have also conducted initial tests on the split gate second generation single devices, in

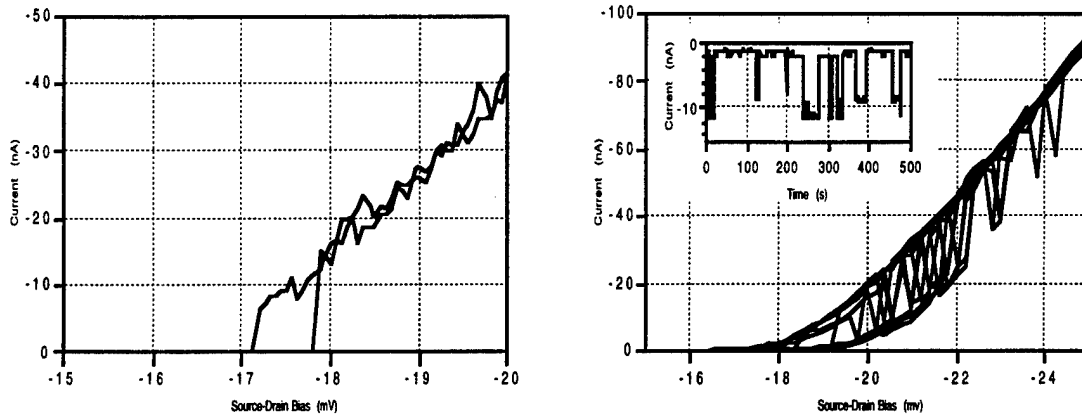


Figure 6: (a) Bistable hysteresis in a CVD point contact at 4.2 K (b) Multi-stable switching and associated multi-level random telegraph signal (inset) in another CVD point contact at 4.2 K.

which the quantum barriers are defined with electrostatic depletion gates instead of wet chemical etching. Well resolved coulomb blockade measurements (Fig. 7) demonstrate that these devices are performing correctly. Future measurements will characterize and compare the switching behavior in this very different architecture to the etched device results.

2.4 Discussion of the Results

The most significant result in the single device investigation has been the characterization of the hysteresis as a basic conduction bistability with a random telegraph signal (RTS). This result has been confirmed in the high bias regime by other groups in an offset split gate [Smith] and a deeply etched lateral barrier [Pilling]. Random telegraph signals have been observed in quantum devices near equilibrium [Dekker][Timp][Sakamoto] and have been attributed to the fluctuations of a single or small number of nearby impurities. Many of our results are consistent with this interpretation, however the exponential dependence of the bistable state lifetimes t_{high} , t_{low} as a function of source drain bias has not been measured before, and remains difficult to interpret

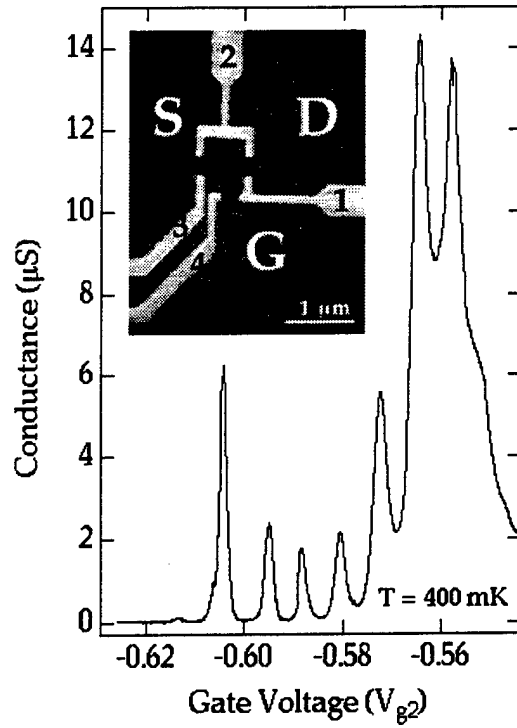


Figure 7: Coulomb blockade oscillations in a three lead dot. The inset shows the SEM picture of the device. Top gates are numbered from "1" to "4". "G", "D" and "S" denote the semi-infinite leads can be used interchangeably as "Source", "Drain" and "Leakage Channel". The result shown here is obtained by varying the voltage on gate "2", and keeping the others constant.

within the impurity model. The very strong control effected by the back gate voltage on switching times is likewise unexplained.

The multi-stability displayed in some of the CVD etched devices (Figure 6) is more typical of fluctuations due to impurities. Yet in this case as before there is an exponential bias dependence of lifetimes, and indeed under controlled circumstances an evolution from bistable on-off switching to multi-stable on-on transitions.

Voltage dependent random telegraph signals have been observed in submicron MOSFET inversion layers and $4\mu\text{m}$ diameter resonant tunnel diodes [Ralls][Ng]. In each case the dependence of the RTS is attributed to the physical position of a switching impurity and it's bias defined energy with respect to a local Fermi level. In our devices the voltage dependence scale is much smaller - the state lifetimes can vary by two orders over only $500\text{ }\mu\text{V}$ of applied bias, inconsistent with the above explanation.

3. Conclusions and Future Work

The cause of the conduction instability remains unclear. Strong qualitative similarities to impurity switching results are contradicted by the exponential voltage dependencies of the state lifetimes. However, we have already been able to demonstrate remarkable control over the character of the instability as it manifests in the I-V relation using both front and back gate potentials. Further probing of this control should lead to a physical explanation of the switching and hysteresis.

We will continue with a series of measurements characterizing the transition from the well understood equilibrium regime to our high bias non-equilibrium situation. Quantitative dependencies of the state lifetimes as a function of gate voltages, applied bias and temperature across this transition are required. Similar measurements on our split gate devices will quantify the relevance of the surfaces and associated imperfections in the etched devices, and direct future fabrication towards the most robust architecture.

With these results in hand, we will return to the performance of the single device arrays, densely packing the point contacts and quantum dots into 1D and 2D arrays. Single and coupled device behavior can then be separated and accurately characterized. This knowledge will form the design framework of future single electron architectures in this regime.

4. References

- [Beenakker] C. W. J. Beenakker *et al.*, *Phys. Rev. B* 44, 1646 (1991)
- [Dekker] C. Dekker *et al.*, *Phys. Rev. Lett.* 66, 2148 (1991)
- [Duruöz] C. I. Duruöz *et al.*, *Phys. Rev. Lett.* 74, 3237 (1995)
- [Harris] J. S. Harris Jr. *et al.*, *JSEP Annual Report* (1994-1995)
- [Kouwenhoven] L. P. Kouwenhoven *et al.*, *J. Phys. B - Cond. Matt.* 85, 367 (1991)
- [Ng] S.-H. Ng *et al.*, *Appl. Phys. Lett.* 62, 2262 (1993)
- [Pilling] G. Pilling *et al.*, *Proceedings EP2DS XI* 347 (1995)
- [Ralls] K. S. Ralls *et al.*, *Phys. Rev. Lett.* 52, 228 (1984)
- [Sakamoto] T. Sakamoto *et al.*, *Appl. Phys. Lett.* 67, 2220 (1995)
- [Smith] J. C. Smith *et al.*, *Proceedings EP2DS XI* 351 (1995)
- [Timp] G. Timp *et al.*, *Phys. Rev. B* 42, 9259 (1990)

5. JSEP Supported Publications

1. C. I. Duruöz, R. M. Clarke, C. M. Marcus and J. S. Harris Jr., "Conductance Threshold, Switching and Hysteresis in Quantum Dot Arrays," *Phys. Rev. Lett.* 74, 3237 (1995).

2. C. I. Duruöz, D. R. Stewart, C. M. Marcus and J. S. Harris Jr., "Switching and Hysteresis in Quantum Dot Arrays," *Proceedings EP2DS XI* 349 (1995).
3. G. Pilling, D. H. Cobden, P. L. McEuen, C. I. Duruöz and J. S. Harris Jr., "Intrinsic Bistability in Nonlinear Transport Through a Submicron Lateral Barrier," *Proceedings EP2DS XI* 347 (1995).
4. G. S. Solomon, C. I. Duröz, C.M. Marcus and J. S. Harris, Jr., "Growth Induced and Patterned 0-Dimensional Quantum Dot Structures" in *Low Dimensional Structures Prepared by Epitaxial Growth or Regrowth on Patterned Substrates*, ed. by K. Eberl et al., NATO ASI Series E, Applied Sciences 298.

6. JSEP Supported Ph. D. Thesis

C. I. Duröz, "Low Temperature Transport in Quantum Dot Arrays", Ph. D. Thesis, Stanford University, March, 1996.

UNIT: 2

**TITLE: Patterned Thin Film Media for
High Density Magnetic Recording**

SENIOR INVESTIGATOR: R. F. W. Pease

RESEARCH STUDENT: R. M. H. New

Background

In conventional hard-disk magnetic recording systems, the signal to noise ratio is often limited by "transition" noise which occurs due to the irregular zig-zag domain walls between adjacent recorded bits [Tong]. In order to address this problem, we are studying recording media composed of large arrays of submicron lithographically defined single-domain magnetic islands. It is known both from theoretical arguments and from experiments that sufficiently small magnetic particles are uniformly magnetized and contain no domain walls. If a single-domain particle of this type has a single uniaxial easy axis of magnetization then it will have only two possible magnetization states and will be ideal for storage of a single bit of information. A magnetic recording medium consisting of an array of equally spaced and uniformly shaped single-domain islands with predictably oriented easy axes could serve as a virtually noise-free alternative to the unpatterned magnetic thin films used in conventional hard disk systems. The ultimate theoretical storage density for such a system would be limited only by the spontaneous thermal switching of bits, a problem that would occur only for particles one hundred angstroms in diameter or less.

In a previous contract period we developed a procedure for patterning polycrystalline magnetic thin films using direct-write electron beam lithography and a multi-step masking and milling process [New (a)]. We used this procedure to define large arrays of $0.15\mu\text{m}$ by $0.2\mu\text{m}$ cobalt islands and studied the physical properties of these islands using atomic force, scanning electron and transmission electron microscopy. The magnetic properties were examined with both magnetic force microscopy and bulk hysteresis loop measurement techniques [New (b)].

For those initial experiments we patterned magnetic islands out of a 200-Å-thick polycrystalline cobalt film. Our results indicated that the transition from the multidomain to single domain state occurs at an island diameter of roughly $0.2\mu\text{m}$. The magnetic force microscopy images of these islands showed that these islands were not single domain. However, smaller islands, roughly $0.15\mu\text{m}$ by $0.2\mu\text{m}$ in size, were almost all single domain. Transmission

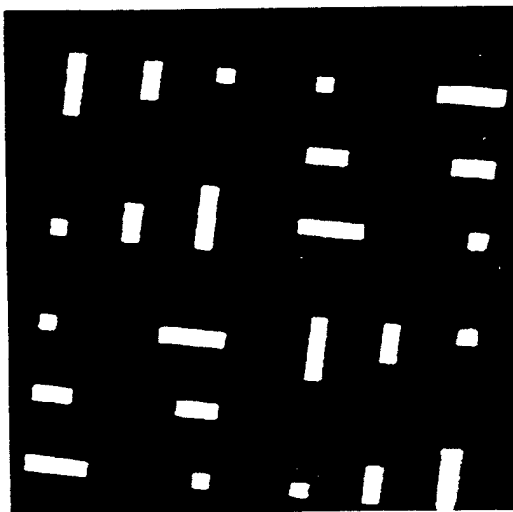
electron microscopy images of the patterned polycrystalline islands indicated that there were roughly 200 cobalt grains per island, each of which has an easy axis of magnetization randomly oriented in the plane of the film. For islands with only a few hundred grains or less, the magnetocrystalline anisotropies of the individual grains may not completely average out and the net magnetocrystalline anisotropy may be larger than the shape anisotropy for some island geometries. Our calculations indicated that for the island geometries we are using, there is a significant probability that the net easy axis may be misaligned with the long axis of the island [New (c)], and our initial experiments confirmed this. Such unpredictably oriented easy axes would cause problems in a single-bit-per-island recording scheme.

One problem with polycrystalline magnetic recording films, either patterned or unpatterned, is that the fundamental unit of magnetization (typically a single grain or grain cluster of 100 to 500 Å in diameter) is not much smaller than the size of a single recorded bit. For a state of the art 1Gbit/in² recording system, there may be only a hundred grain clusters or less per bit. Because the medium is so coarsely discretized, conventional magnetic recording systems suffer from increasing signal to noise problems as recording densities are increased. Medium noise is already the most important component of noise in recording systems that use magnetoresistive readback heads.

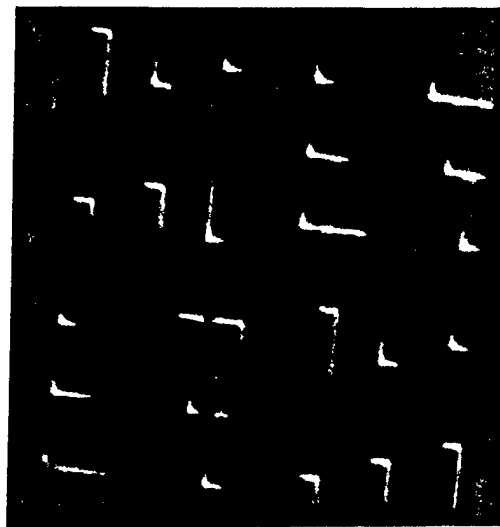
One solution to this problem is to switch to a recording medium which is homogeneous over the size range of a single recorded bit. Sputtered single crystal films would provide a more controllable and predictable magnetic behavior when examined at this size range, and patterned islands of single crystal material would not suffer from the problem of randomly oriented easy axes discussed above. One of the major advantages of the patterning technique that we have developed is that, unlike a lift-off process, it can be used to pattern single crystal thin films. In preparation for future experiments we have sputter-deposited single crystal iron films on sapphire substrates and measured their magnetic and structural properties. These films show good epitaxial quality and have a predictably oriented uniaxial anisotropy as required.

Progress during the current Year

During the last year we successfully patterned such films into islands with lateral dimensions ranging from 100nm to several microns, using techniques not dissimilar to those used for patterning the cobalt. The resulting islands were examined using a variety of techniques including SEM, AFM, alternating gradient magnetometer (Princeton Measurements Corp.) and vibrating sample magnetometer (Kobe Steel). In all cases the islands were single crystal. This included islands up to several microns long and with high aspect ratios with the long axis at a



(a)



(b)

Figure 1: Topographical (a) and magnetic (b) images of single-crystal iron islands with a smallest feature size of 1.0 microns. Even the largest islands (up to 1.0 by 3.0 microns) are single domain.

variety of orientations (Fig. 1). All of the islands had their easy axis aligned with the surface anisotropy of the film thus we were able to control simultaneously the crystallographic orientation and the shapes of the islands and hence compare the relative strengths of magnetocrystalline and shape anisotropy of these islands. In addition to this scientific advantage it is important technically to have the easy axis across the long direction of the islands because it allows for more efficient coupling of the medium to the read/write head. This can be accomplished when crystal orientation is the dominant factor governing easy axis direction; if shape anisotropy dominated then the easy axis would be along the long direction.

We have also, in collaboration with University of Maryland, been able to examine these islands in a magnetic contrast AFM while applying an external field and observe the external field necessary to switch each island (Fig. 2). By examining a large population of islands under these conditions we were able to predict quantitatively the shape of the magnetization loop of the complete sample. This prediction turned out to be accurate thus confirming our model of the contrast mechanism for the magnetic contrast AFM.

We have developed a preliminary model for the reduction in transition jitter that might be expected if we employed a patterned recording medium. Current recording systems suffer from a transition jitter described as a standard deviation σ_t of the transition point of about 5nm for a track width of 3 μ m. Narrower track widths will show worse jitter because of the lower number of grains being averaged over. With a patterned medium we might expect the jitter to be consistent with the edge roughness when averaged over the length of the island and this could well be a factor of 4 lower. More comprehensive models would have to consider the particular form of read head used and the signal extraction algorithms used.

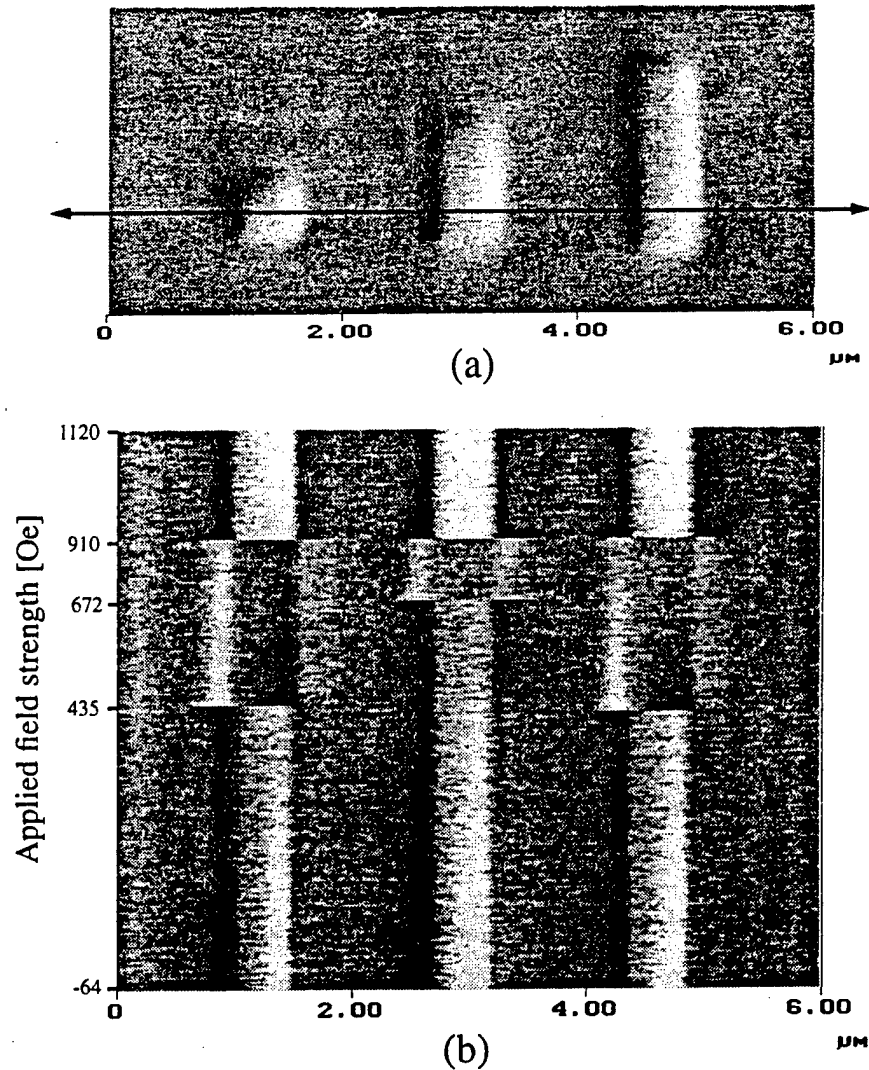


Figure 2: Figure (a) shows a magnetic force image of three islands, each with a width of $0.5\mu\text{m}$, but with different lengths. In Figure (b), the vertical scan direction has been turned off, and the magnetic tip moves back and forth along a horizontal line (\longleftrightarrow) as the externally applied field from the electromagnet is ramped up. The field is applied within five degrees of the easy axis of magnetization. Both the electromagnet and the islands are initially in a saturation remanent state (at the bottom of (b)). As the field is increased, the left and right islands are the first to reverse their magnetizations (at 435 Oe) followed by the center island (at 672 Oe). At 910 Oe, the magnetization of the tip reverses so as to align itself with the applied field.

During the reporting period the student, Richard M. H. New, completed his PhD. requirements and graduated and is now at the IBM Almaden Research Center San Jose CA. His dissertation, "Patterned Media for High Density Recording", was approved in September 1995 and copies are available.

References

- [New (a)] R. M. H. New, R. F. W. Pease, R. L. White, *J. Vac. Sci. Technol. B*, **6**, 3196, Nov/Dec 1994.
- [New (b)] R. M. H. New, R. F. W. Pease, R. L. White, *J. Vac. Sci. Technol. A*, May/June 1995.
- [New (c)] R. M. H. New, R. F. W. Pease, R. L. White, submitted to IEEE International Magnetics Conference, April 1995.
- [Tong] H. C. Tong, R. Ferrier, P. Chang, J. Tzeng and K. L. Parker, *IEEE Trans. Mag.*, **20**, 5, 1831 (1984).

JSEP Supported Publications

1. "Magnetic force microscopy of single-domain single-crystal iron particles with uniaxial surface anisotropy," R. M. H. New, R. F. W. Pease, R. L. White, R. M. Osgood, K. Babcock, to be published in the *Proceedings of the 40th Annual Conference on Magnetism and Magnetic Materials (J. Appl. Phys.)* held in Philadelphia, Nov. 1995.
2. "Lithographically patterned single domain cobalt islands for high density magnetic recording," R. M. H. New, R. F. W. Pease, R. L. White, to be published in the *Proceedings of the 6th International Conference on Magnetic Recording Media (J. Magn. Mater.)*, held in Oxford, England, July 1995.
3. "Effect of magnetocrystalline anisotropy in single-domain polycrystalline cobalt islands," *IEEE Trans. Mag.*, MAG-31, p. 3805, Nov. 1995.

JSEP Support Thesis

"Patterned Media for High Density Magnetic Recording," R. M. H. New, Ph.D. Thesis, Stanford University, September, 1995.

UNIT: 3

**TITLE: Investigation of a Metal Source and Drain
Field Emission Transistor**

PRINCIPAL INVESTIGATOR: C. R. Helms

GRADUATE STUDENT: J. P. Snyder

Background

Metal source and drain Metal-Oxide-Semiconductor-Field-Effect-Transistors (MOSFETs) have been shown to have several key advantages over their conventional (doped source and drain) counterparts including ease of fabrication and unconditional immunity to parasitic bipolar and latch-up effects. They were first investigated in the late 1960s [Lepselter], and were thought to have certain advantages over their conventional (diffused source and drain) counterparts including a simplified process, the ability to make very shallow source and drain regions, low source and drain sheet resistance, and complete immunity to latch-up and parasitic bipolar effects. They proved to be poor performers however when compared to a similarly sized conventional MOSFET. The lower drive current in the 'on' state was attributed to the presence of a finite 'gap' between the edge of the poly gate and the edge of the platinum silicide (PtSi) source metal. The much higher leakage currents in the 'off' state originate at the drain end of the device, where electric fields cause the thermally assisted field emission of electrons from the drain into the silicon [Lepselter] [Oh] [Koenekke] [Sugino] [Tsui].

Until recently, the low temperature characteristics of these devices have not been investigated. The only exception to this is a 1968 paper [Lepselter] in which 77 K I-V curves are shown and briefly discussed. Their device was fabricated with a non-self aligned, chemical vapor deposition (CVD) gate oxide process. The data shows a significant *decrease* in current drive at 77 K compared to room temperature.

Since 1993, several papers [Tucker] [Hareland] have reported on simulations on these and similar devices, and have shown acceptable drive current and short channel effects in devices with channel lengths down to 0.025 μm . The scalability of these metal source and drain devices is particularly impressive at low temperatures (77 K), as described by [Tucker]. It seems possible in

light of these recent studies to build a metal source and drain device that has all the advantages previously mentioned, as well as superior scalability to well below $0.1\ \mu\text{m}$ and free of the low drive and high leakage current problems. The only requirement is low temperature operation.

Progress during the current Year

We report the first detailed experimental investigation of the low temperature, field emission characteristics of PtSi source and drain MOSFETs. I-V curves have been measured at

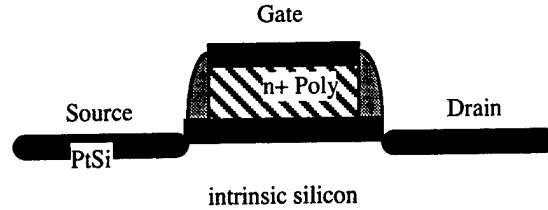


Figure 1: Schematic Diagram of the Device.

various temperatures down to 4.2 K and for channel lengths down to $1\ \mu\text{m}$. Device fabrication has been optimized so that it is free from the 'gap' at the poly edge described earlier. As will be discussed, we observe a definite transition in the current flow mechanism of the device, from thermal to field emission, as the temperature is reduced below 100 K. In this low temperature 'field emission mode', the drive current when the device is 'on' is comparable to that of a conventional MOSFET, and short channel effects are not observable down to $1\ \mu\text{m}$, despite the fact that the substrate is nominally undoped. The schematic diagram of the device is shown in Fig. 1.

The band diagrams of Fig. 2 demonstrate the operating principle of the device described in figure 1 at an intermediate temperature ($\sim 150\ \text{K}$) such that the various current flow mechanisms are observable. The band diagrams are drawn along a line from source to drain, just underneath the gate oxide, and show the Fermi levels of the source and drain PtSi, as well as the conduction and valence bands of the silicon substrate.

In Fig. 2(a), when the device is in its 'off' state with bias applied only to the drain, hole leakage current enters the channel by thermal emission over the sum of the 0.2 eV Schottky barrier and an electrostatic barrier present because of the difference in workfunction between the n+ polygate and the PtSi source. In this domain of gate voltage (V_g), the thermal emission regime, holes flow by diffusion from source to drain and the silicon bands in the channel are flat. Changing the gate voltage simply changes the amount of hole thermal emission current entering the

channel, as is seen in the 'thermal emission characteristic' drawn in the plot of source current (I_s)

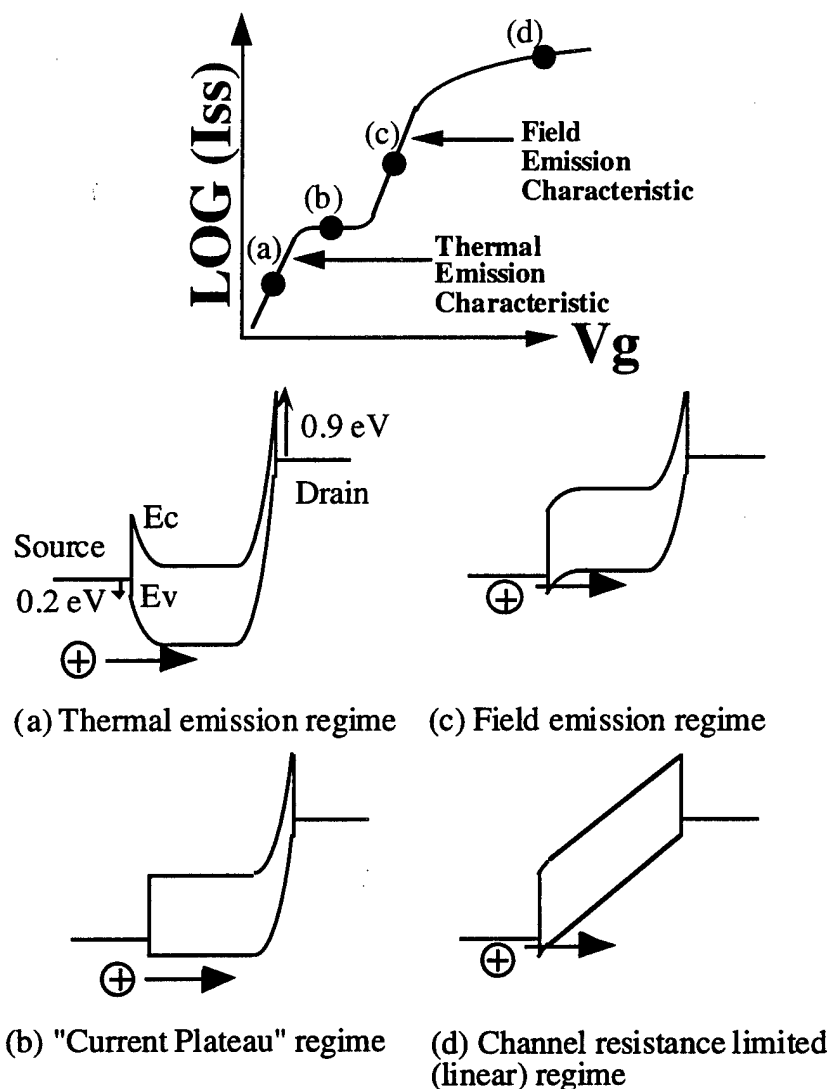


Figure 2. A band diagram description of the different current flow regimes seen in a typical source current vs. gate voltage plot. (a) Thermal emission regime (b) "current plateau" regime (c) field emission regime and (d) channel resistance limited regime.

vs. V_g . There is also the possibility of electrons being field emitted from the drain because of the high electric fields there, but this component of current does not show up in our measurements of source current and will not be discussed in this report.

Eventually, with increasingly negative gate bias, only the fixed Schottky part of the barrier to holes remains and the current is limited by thermal emission over this barrier [Fig. 2(b)]. In this

'current plateau' regime further increases in the magnitude of the gate voltage cease to have an exponential effect on I_s . The hole current is, for the most part, dependent only on the temperature and the barrier height (~ 0.2 eV), as is drawn in the topmost plot.

With high enough gate bias, holes eventually can be made to tunnel through the Schottky barrier and I_s once again begins to increase in an exponential fashion, this time along a 'field emission characteristic' [Fig. 2(c)]. The current is not yet large enough to give the silicon bands in the channel appreciable slope, which is to say that the current is still field emission limited and still travels by diffusion from source to drain, and is not yet channel resistance limited.

Finally I_s becomes large enough that the channel resistance begins to dominate and the holes travel by drift [Fig. 2(d)]. In this regime of V_g the current drive of the device is similar to that of a conventional MOSFET as the Schottky barrier has been rendered all but transparent to the flow of holes.

Drain curves (I_s vs. drain voltage (V_d)) and gate curves (I_s vs. V_g) were measured with a computer controlled HP 4140B DC voltage source/pA meter. A Lakeshore cryogenic probe station was used to perform measurements down to 4.2 K.

Figure 3(b) shows the experimental gate curves of the device described in Figs. 1 and 2 with width=length=2 μm . Here the thermal emission, plateau, field emission and channel resistance limited regimes are clearly seen, especially for the 200 K curve. As was mentioned previously, the plateau current is solely a function of temperature and barrier height and this dependence is observable. The plateau current drops exponentially with temperature, so that for temperatures less than about 100 K, all significant current flow (> 0.1 pA) occurs by the process of field emission and the device is being operated in the 'field emission mode'. It can be seen that this field emission characteristic is largely independent of temperature. Because n+ poly is used for the gate material, V_g must be brought to about -2 Volts before significant current begins to flow. Referring back to Fig. 2, this implies that even the band diagram in Fig. 2(c) could be used as an effective 'off' state. This could be realized for example, if p+ poly were used for the gate. The Schottky barrier alone is responsible for preventing the flow of current into the channel and thus it is clear why substrate doping is not required.

It is also possible to back out the effective PtSi - Si barrier height to holes from the thermal emission formula $I = AA \cdot T^2 \text{Exp}(q\phi_b/kT/(kT/q))$ using the plateau currents and corresponding

temperatures. This formula gives a barrier of ~ 0.195 eV, in very good agreement with published

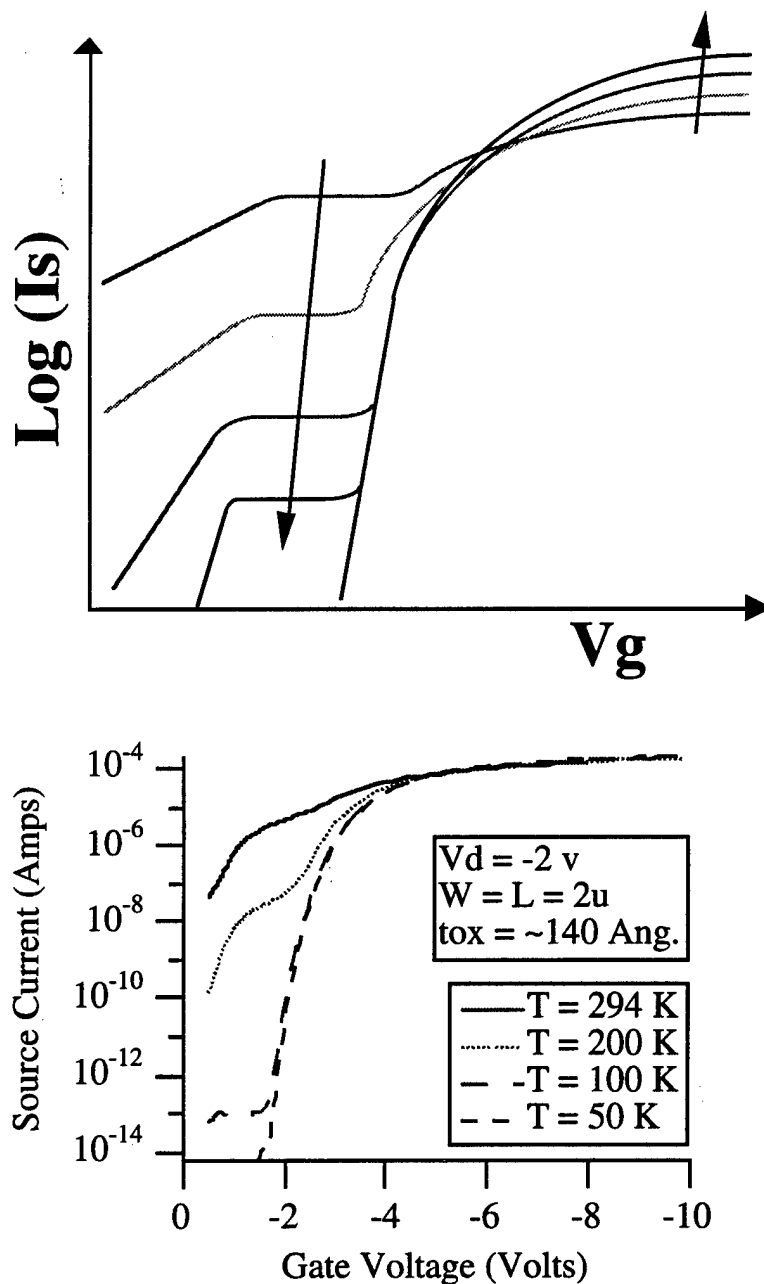


Figure 3. Variation with temperature (a) qualitative example showing the major effects of temperature variation on the gate curves of a PtSi source and drain MOSFET. The arrows point in the direction of decreasing temperature. (b) Actual measured data of a device described in (a).

barrier heights of the PtSi - Si system [Mooney] [Weeks].

During the last year also developed a full 2-D Poisson solver which is integrated with first principles tunneling calculations in order to theoretically examine the effects of device geometry (tip sharpness, channel length, and gate oxide thickness) and materials and system parameters (Schottky barrier height and temperature) on the hole and electron field emission characteristics. The subthreshold slopes of these characteristics were found to decrease monotonically with gate oxide thickness with no theoretical limit. This is in contrast to the theoretical limit, defined by temperature, that exists for the subthreshold region of a conventional device. Subthreshold current levels were also found to be generally smaller than those of conventional devices by several orders of magnitude. Shallow source/drain junctions with sharp tips were found to be optimal in terms of promoting hole field emission drive currents and controlling Drain-Induced-Barrier-Thinning (DIBT) hole leakage currents. Low barrier heights (for good drive currents) and low temperatures (for low leakage over the low barrier) were also found to be optimal.

Possible Future Directions

These devices will be investigated further. Shallower junction, p+ poly, no gap devices (unlike the ones studied in this dissertation) will be investigated especially with regard to drive current and electron leakage current. NMOS devices can be built as long as metal-silicon Schottky diodes with low barriers to electrons can be found. Rare-earth silicides are potential candidates for this application. Finally, full 2-D modeling of these field emission devices with integrated tunneling and hot-carrier models will be used to further explore the 'virtual source voltage' phenomena described in Chapter 8 of J. P. Snyder's Ph.D. thesis, and to determine the effects of this phenomena on device long term reliability.

References

- [Hareland] S. A. Hareland, A. F. Tasch, C. M. Maziar, *Electronics Letters*, **29**, 1894 (1993).
- [Hareland] S. A. Hareland, A. F. Tasch, C. M. Maziar, *Proceedings of the 21st International Symposium on Compound Semiconductors*, September 18-22, San Diego, CA (1994).
- [Koenekke] C. J. Koenekke, S. M. Sze, R. M. Levin, E. Kinsbron, 1981 IEDM, 367.
- [Lepselter] M. P. Lepselter, S. M. Sze, *Proceedings of the IEEE*, 1400 (1968).
- [Mooney] J. M. Mooney, J. Silverman, M. M. Weeks, *SPIE, Infrared Sensors and Sensor Fusion*, **782**, 99 (1987).
- [Oh] C. S. Oh, Y. H. Koh, C. K. Kim, 1984 IEDM, 609.
- [Sugino] M. Sugino, L.A. Akers, M.E. Rebeschini, 1982 IEDM, 462.
- [Tsui] B. Tsui, M. Chen, *J. Electrochem. Soc.*, **136**, 1456 (1989).

- [Tucker] J. R. Tucker, C. Wang, J. W. Lyding, T. C. Shen, G. C. Abeln, *1994 SSDM*, 322.
- [Tucker] J.R. Tucker, C. Wang, P.S. Carney, *Appl. Phys. Lett.*, **65**, 618 (1994).
- [Weeks] M. M. Weeks, P. W. Pellegrini, SPIE, Test and Evaluation of Infrared Detectors and Arrays, **1108**, 31 (1989).

JSEP Supported Publications

- J. P. Snyder and C. R. Helms, Y. Nishi, "Experimental investigation of a PtSi source and drain field emission transistor," *App.Phys.Lett.* **67**(10), 4 September 1995.

UNIT: 4

TITLE: On-Chip Thin Film Solid State Micro-Battery

PRINCIPAL INVESTIGATOR: S. S. Wong

GRADUATE STUDENT: J. Leung

Scientific Objectives

The objectives of this work are to develop the fabrication technology and characterize the performance of thin film solid state micro-batteries that are suitable for monolithic integration with semiconductor devices.

Summary of Research

In the last phase of this research, we successfully fabricated and tested solid-state micro-batteries on a silicon substrate. Charging and discharging up to 1000 cycles were demonstrated. In addition, we identified a suitable diffusion barrier, PECVD oxynitride, that can prevent the lithium ions from penetrating down to the devices in the substrate.

Encouraged by the experimental results, we have designed various circuits to be integrated with the micro-batteries. These circuits include charging and discharging control units, various temperature sensors, and voltage controlled oscillators. Figure 1 shows a basic charging circuit. Both constant and pulse current charging are amongst the various charging modes to be studied. We have also taken into account the possibility of circuit leakage. Simulation indicates that the leakage current should be less than 3 pA when the micro-battery is fully charged to above 2V. Both internal resistance and capacitance of the battery are included in the modeling. The basic op amp is capable of operating in the high MHz range while delivering 200 uA. In addition, various temperature sensors will be placed on chip. The close proximity of these sensors to the battery will yield accurate thermal profile which will be important in the design of large integrated circuits.

We will also incorporate a voltage controlled oscillator (VCO). The VCO operates off the on-chip battery and can be programmed with an external signal. Simulation has indicated that the VCO can run up to 1 MHz. The excellent voltage stability of the micro-battery will be ideal for applications such as frequency synthesizer which requires VCO with an excellent power supply rejection ratio. Noise at VCO is often difficult to be filtered out in a frequency synthesizer, and is the factor that determines channel spacing criteria [Mannasiwitsch].

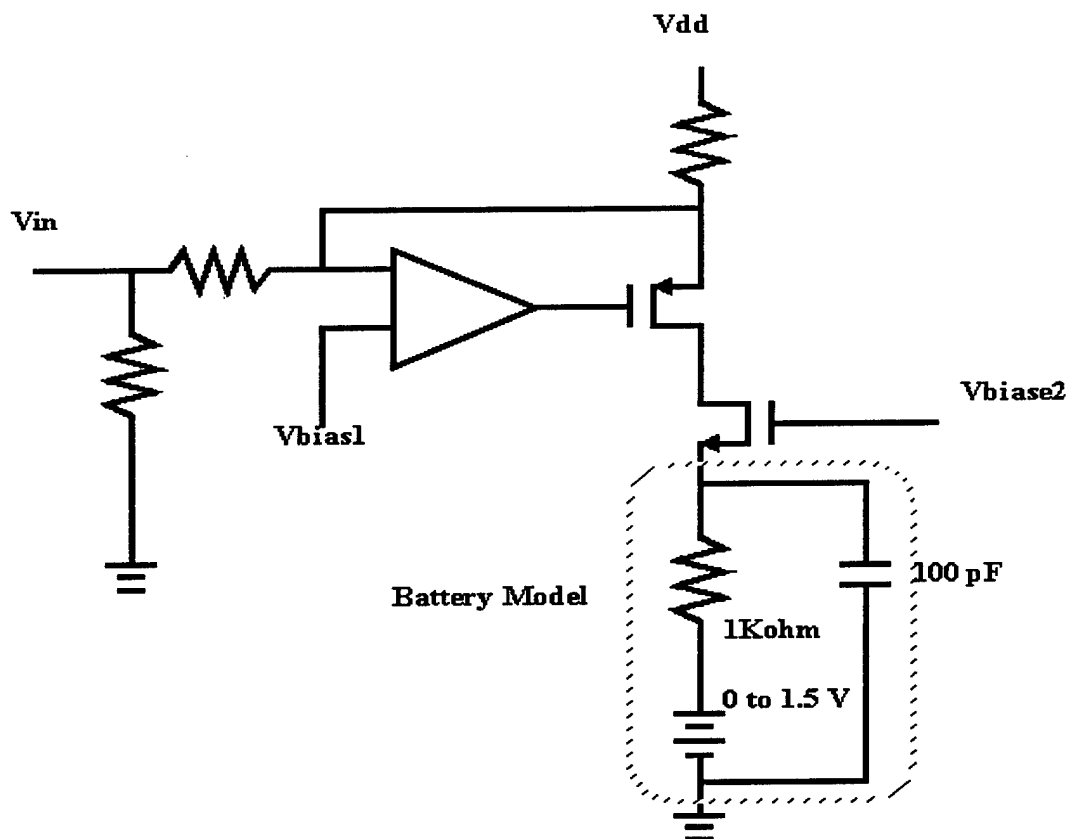


Figure 1. Basic charging circuit.

Figure 2 illustrates the sequence for monolithic integration. The circuits will be first fabricated with a conventional CMOS technology. Afterwards, a layer of silicon oxynitride passivation layer will be deposited using plasma enhanced chemical vapor deposition (PECVD). Lastly, the various layers for the lithium battery will be sputtered on.

The circuits will be fabricate on 4-inch wafers in a 2 μm CMOS technology. Individual die size is limited to about 8 μm by 8 μm . Ten micro-batteries will be sputtered on each wafer. Each battery will be about 1 cm by 1 cm and with a charge capacity of about 1 Coulomb. An overview of the wafer is shown in Fig. 3.

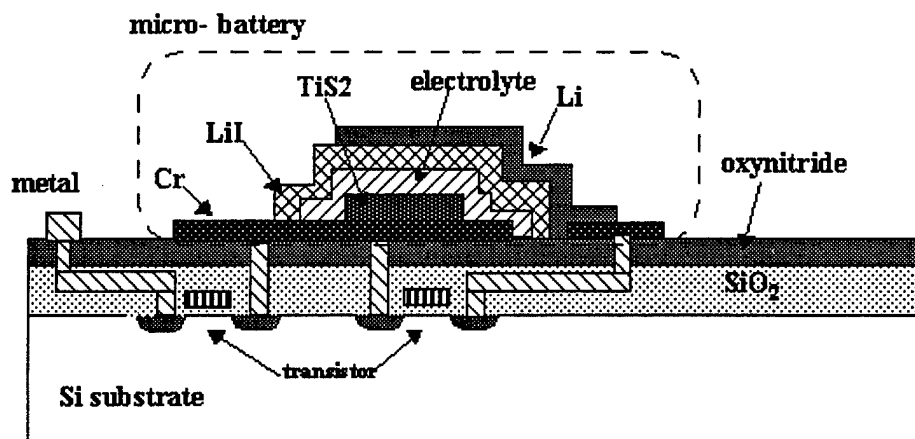


Figure 2. A cross section of the integrated micro-battery.

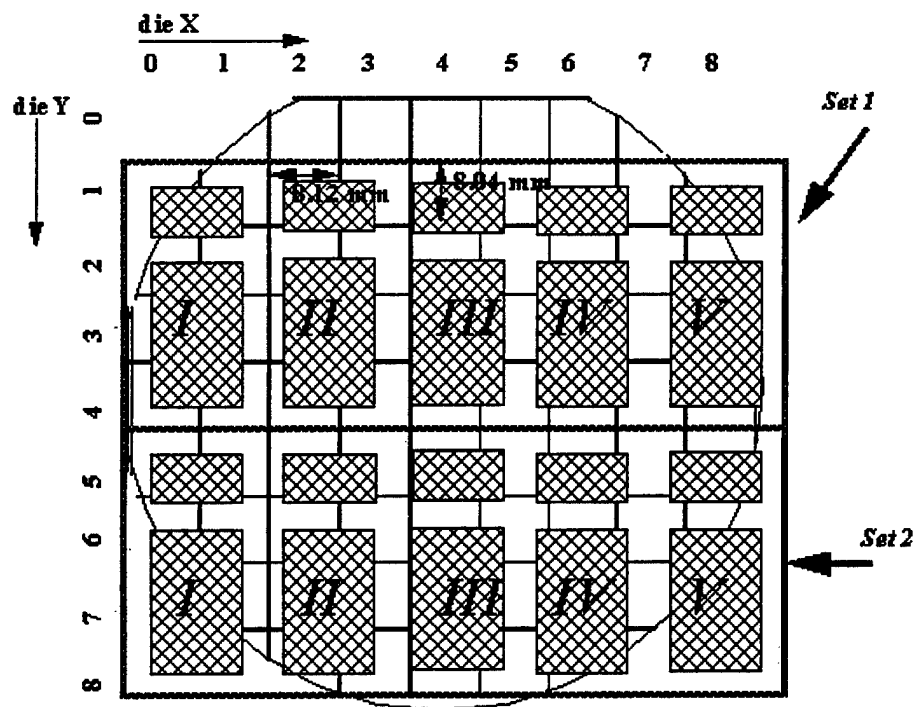


Figure 3. Placement of micro-batteries on a four-inch wafer.

The fabrication will commence shortly. We aim at not only demonstrating the feasibility of such a full integration, but also performing on-chip characterization.

References

[Mannasiwitsch] V. Manassewitsch, Frequency Synthesizer, Wiley, 3rd Edition.

UNIT: 5

TITLE: CVD Epitaxial Germanium *n*-Channel FETs Formed on Si Substrates using Strain-relief Layers

PRINCIPAL INVESTIGATOR: K. Saraswat

GRADUATE STUDENT: D. Connelly

Abstract

N-channel field effect transistors are fabricated in strained and unstrained Ge grown via graded-alloy strain reduction on (001) silicon substrates. Applications of Ge device integration with silicon substrates are discussed. Blanket graded-alloy epitaxy is compared with other strain reduction techniques. The effect of strain on the Ge conduction band structure and hence on electron transport in the *x*-*y* plane is examined.

Objectives

The following are the primary objectives of this project::

- To fabricate *n*-type Ge-channel MOSFETs on a Si substrate.
- To investigate the effect of different degrees of compressive strain on the electron transport properties in germanium inversion layers.
- To compare different schemes for the formation of strain-relief structure formation including blanket graded epitaxy, selective graded epitaxy, and graded epitaxy on ultra-thin silicon-on-insulator.
- To assess the utility of high-germanium content *n*-channel MODFETs in high-speed transistor applications.

Prior Art

The development of strained layer epitaxy of GeSi alloys on silicon substrates sparked interest in the development of heterostructure devices using silicon-based technology. Much of the work can be placed in one of two categories, vertical heterostructure bipolar transistors (see for example [King]), in which the primary interest is the band-gap difference between the base alloy and the emitter alloy, and confined-carrier field-effect devices (see for example [Pearsall86] and [Daembkes]) in which the parameter of interest is the conduction band offset (for *n*-channel devices) or the valence band offset (for *p*-channel devices).

The biaxial compressive strain formed when GeSi with non-zero x is deposited on silicon enhances the natural positive valence band offset of the GeSi relative to silicon [Pearsall89]. Representative is the work by the group from UCLA [Nayak] in which a 10 nm "undoped" unstrained silicon layer is deposited on an n -type Si substrate. An undoped 15 nm strained $\text{Ge}_{0.20}\text{Si}_{0.80}$ layer is then deposited to form the channel region. It is capped with a 10 nm silicon layer. A 5 nm SiO_2 layer is then thermally grown to form the insulator, consuming some of the underlying silicon. The structure is capped with a polycrystalline silicon gate electrode. The estimated 0.15 eV valence band discontinuity confines most of the holes to the $\text{Ge}_{0.20}\text{Si}_{0.80}$ layer for the initial portion of the superthreshold gate bias regime. The Princeton group [Garone] fabricated a similar structure with a 10 nm $\text{Ge}_{0.40}\text{Si}_{0.60}$ well capped by a 7.5 nm silicon spacer and a 10 nm gate oxide with an aluminum gate electrode.

For electron-confinement structures things are more complicated. For unstrained material of low-to-moderate germanium concentrations the conduction band consists of six degenerate ellipsoids aligned along the x , y , and z axes in momentum space. For material under [001]-directed strain the degeneracy is broken --- the z -directed ellipsoid is either raised or lowered in energy relative to the x and y valleys. Since the z valley exposes its carriers during conduction in the x - y plane to only the light transverse effective mass it is preferable to raise the energy and thereby reduce the carrier population of the x and y ellipsoids. This is done by depositing the channel material in biaxial tension.

Reported work to date has been on structures utilizing strained silicon as the channel material. Representative is the work from Stanford [Welser]. They had two forms of their device. One started with a relaxed (001) $\text{Ge}_{0.30}\text{Si}_{0.70}$ surface on which was grown a strained silicon layer. Subsequent oxidation of the silicon resulted in a 12.8 nm gate oxide over a 4.6 nm strained silicon channel well. The other used a $\text{Ge}_{0.29}\text{Si}_{0.71}$ surface on which was grown an 8.0 nm strained silicon layer covered with a 7.2 nm $\text{Ge}_{0.29}\text{Si}_{0.71}$ spacer and a "sacrificial" strained silicon cap. Thermal oxidation to form the 12.8 nm gate oxide fully consumed the cap. Thus the former devices had surface channels while the latter had buried channels for moderate superthreshold biases.

IBM has published results of Schottky-gated n -channel structures using both Molecular Beam Epitaxy and Ultra High Vacuum Chemical Vapor Deposition (UHVCVD) [Ismail] [Wang]. They used a starting relaxed surface with a 30% germanium content. Their channel was formed in a 10.6 nm strained silicon layer.

The key difficulty in the formation of these structures is the preparation of the initial surface. Ideally if a surface of a given alloy composition is needed a wafer of that composition should be used. Unfortunately wafers of arbitrary germanium content are not available --- silicon wafers are widely available and germanium wafers are available at considerably higher cost. A solution is to deposit a relaxed "buffer layer" in which threading dislocations are isolated below the surface to translate the surface composition to the desired value from what is on hand. Leaders in this technique include AT&T with Molecular Beam Epitaxy and IBM with CVD and MBE. All reference cases described here begin with (001) silicon wafers.

In reference [Fitzgerald] AT&T reports the results of linearly grading the germanium content from zero up to 53% using MBE at 900 C. The alloyed germanium content is ramped at 0.1 per micrometer. The high temperature is used to prevent the accumulation of stress in their films before relaxation events, increasing the number of threading dislocations and canceling the benefits of compositional grading. They fabricate ungated electron confinement structures [Xie] with good results. Other workers [Schaffler] showed that increasing the gradient to 0.45 per micrometer and decreasing the deposition temperature to 750 C can still yield significant advantage over abruptly stepped buffer layers.

IBM has generated strain-relief layers using both continuous grading as per AT&T (see [Legoues91]) and grading in discretized steps [Meyerson] using both MBE and CVD. In reference [Tsang] step-grading from pure silicon to pure germanium with fewer than 0.01 threading dislocations per square micrometer in the top germanium film is reported. The deposition is done using UHVCVD with composition graded in 40 steps at approximately 0.20 per micrometer. The quality of the film is sensitive to the deposition temperature, with 450 C optimal for the pure-germanium portion. This is described in reference [Legoues92].

The Stanford group [Welser] used, for example, a graded layer of germanium composition from 6% to 30% continuously graded over 1.6 μm deposited via CVD using "Limited Reaction Processing" at 750 C. They had difficulty grading beyond 50% germanium starting from pure silicon using their technique, although IBM's positive results show it can be done.

Background

Single-crystal GeSi alloy exhibits a peak valence band energy that increases steadily with increasing Ge content. The energy of the sixfold-degenerate X (used here to signify all six $\langle 100 \rangle$ directions) conduction band valleys is relatively insensitive to the Ge content in unstrained material. Up to approximately 80 atomic-percent Ge these X-valleys have the lowest energy of the

conduction states in the material. At higher Ge concentrations, however, the strong alloy-dependence of the eight-fold degenerate $\langle 111 \rangle$ L -valleys brings them to a lower energy.

Due to the dependence of the valence band energy on alloy content across the material spectrum most unipolar heterostructure devices built in the low-Ge regime have used holes as their carrier. n -type devices have been built, however, exploiting the strain-dependence of the conduction band minimum.

When (001) silicon is deposited pseudomorphically on a thick unstrained crystalline GeSi alloy the silicon is in biaxial tension, with decreased lattice spacing in the growth direction (z) and increased lattice spacing in the two orthogonal directions (x and y). The result is that electrons in the z -valleys ([001] and [00-1]) exhibit a reduced energy relative to those in unstrained silicon while the x and y valleys see an increase in the energy of their states. (See [Pearsall89] for a good overview of the strain effects on GeSi bands.) The advantages are two-fold. First, since the unstrained GeSi substrate has similar conduction band energies to unstrained silicon, the Si now has a reduced conduction band energy relative to the surrounding material and electron confinement can be achieved. The second advantage is that these valleys exhibit a transverse effective mass lower than their longitudinal effective mass. Since conduction in the channel by z -valley electrons will be characterized by the lower transverse effective mass while electrons in the other four valleys will be subject to a mixture of the longitudinal and transverse effective masses, preferential occupation of the z valleys results in a decrease in net effective mass and a corresponding increase in mobility for appropriate carrier densities. The stress-induced electron confinement for devices in principle works for alloys from zero Ge up to approximately 80 atomic percent Ge. However, work to date has focused on using strained silicon as the channel material.

In Ge-rich material there is therefore available two mechanisms to yield band offsets. If the unstrained starting material is (001) $\text{Ge}_{0.75}\text{Si}_{0.25}$ then application of a strained layer of pure Ge will result in a reduced conduction band energy due to the lower energy of the L -valleys (due to symmetry the effect of the [001] compression on the $\langle 111 \rangle$ L -valleys is small). Growth of a strained $\text{Ge}_{0.50}\text{Si}_{0.50}$ film on the same substrate will result in reduction of the z -valley energies relative to the unstrained material. These offsets could be used in the formation of confined-electron structures.

Of further interest in Ge channel devices is in which valleys the conduction band minimum occurs. As the degree of [001] compression is increased via a lowering of the effective substrate germanium content, the energy reduction of the x and y valleys increases the population of

electrons occupying them until they become the principle repository for channel electrons. The effect of this transition on electron mass and electron scattering is of significant importance.

Of practical interest is the formation of the relaxed buffer layer. Linear grades can be done via different temperature schedules to confine stress-relieving defects below the surface. These grades can be executed either on a blanket wafer or in regions defined in a surface oxide layer. Another option is the formation of a graded buffer layer on ultra-thin silicon-on-insulator, decreasing the energy needed to relax the surface.

Current Work

While there are many interesting possibilities with Ge-on-Si devices, due to the considerable challenges encountered in the optimization of the graded epitaxial process and in the reliable formation of dielectrics on a germanium surface, this project is focusing on two, both currently under fabrication. One is simple Ge-on-Si *n*-channel field effect transistors. These are expected to exhibit conduction-band minima in the *L*-valleys such as those exhibited by bulk germanium, as was discussed in the last section. The second type of device is the strained Ge-channel on strain-reduced GeSi using a germanium atomic fraction of 75%. It is expected that the strain will reduce the energy of *x* and *y*-directed delta-points below the *L*-valleys, yielding a significant and observable difference in in-plane carrier transport.

Strain-relief via graded epitaxy is achieved by grading the composition, pressure, and temperature in the epitaxial reactor. Depositions are done in the Stanford Center for Integrated Systems Applied Semiconductor Materials Epsilon Chemical Vapor Deposition Epitaxial Reactor. The reactor is a multi-lamp-heated single-wafer unit with a graphite susceptor.

Starting wafers are 4-inch 10 ohm-cm boron-doped (001) silicon. These are cleaned via the lab's standard "HF-last" prediffusion clean and immediately placed in the reactor load lock. After at least an hour's nitrogen purge, the wafers are "prepared" by lowering the load-lock elevators into the exchange chamber. The processing of each wafer starts with an atmospheric-pressure hydrogen bake at 1150 C, an *in-situ* H₂+HCl etch at 1150 C, and another 1150 C hydrogen bake.

Germanium grading is achieved by ramping up the germane flow, ramping down the silane flow, and ramping down the temperature continuously during the deposition period. At the end of the grade, a germanium cap of approximately 2 μm is deposited for device formation. Mass flow controller limitations bound the contiguously-available range of germanium fractions at 3% and

98%, however the "discontinuous" jumps from 0 to 3% and 98% to 100% are accommodated without noticeable quality degradation in the film quality.

The key to successful strain relaxation is to maximize the strain reduction achieved via the formation of buried misfit dislocations. These nucleate either homogeneously (thermally) or heterogeneously (due to external factors, such as particles, the wafer edge, etc.). These misfits generally form and propagate in either the [110] or [1-10] direction until either the temperature drops below a kinetic threshold, the edge of the epitaxial region (the wafer edge in the case of blanket epitaxy) is reached, or they scatter towards a wafer surface in the form of threading arms. Since threading arms at the surface can degrade device performance, the distance the misfits are able to travel before scattering should be maximized. A combination of high deposition temperature to drive the propagation kinetics, low deposition rate to give the misfit time to propagate, and low growth rate to maintain an acceptable level of residual strain is thus desirable.

Low deposition rate is accomplished by keeping the partial pressures of silane and germane low. However, the combination of a low deposition rate and a gentle alloy gradient yields long deposition times, a potential practical impediment. High deposition temperature causes other problems. Gas phase nucleation, which causes particulate contamination of the surface and formation of a non-epitaxial film, is activated with temperature. Another practical problem with high deposition temperatures is coating of the chamber wall can occur. Since stopping the deposition in-progress is undesirable, it is important that chamber deposition be kept sufficiently low that quartz transparency is maintained.

The primary tools used for material quality determination, other than device fabrication, have been AFM, TEM, RAMAN spectroscopy, EMP, RBS, and anisotropic etches. AFM is of particular interest, as it can be done nondestructively with rapid turnaround on the full-wafer Park Scientific atomic force microscope in the Stanford Center for Integrated Systems. The strain reduction process results in surface undulations in the material. When grading is done from silicon to pure germanium, the peak slope of these undulations is approximately one degree with a mean spacing between local peaks of order 5 to 10 micrometers. These are the result of the system's attempt to minimize energy -- when the equilibrium mean lattice spacing of an alloy being deposited is greater than the available mean lattice spacing of the exposed alloy surface, the system uses its degree of freedom in the z -direction to increase the mean spacing between deposited atoms. This yields coherent surface undulations in the [110] and [1-10] directions on the surface. For films deposited at sufficiently high temperature, sufficiently shallow alloy gradient, and at sufficiently low deposition rate, these undulations extend for thousands of micrometers. On films deposited under less optimal conditions, these undulations can be quite short, even 10 micrometers or less, at

which point their orientation becomes difficult to determine. Another indicator of poor quality is observed in films deposited with an excessive temperature schedule -- round pits appear in the surface. These are suspected to be due to gas-phase nucleation yielding particulate contamination of the surface and a resulting disruption of "uniform" epitaxial deposition.

Since the source and drain of the FETs are *n*-type, *p*-type doping for the body is needed. The substrate is thus boron doped, and diborane is flowed with the germane during formation of the germanium cap to yield a boron concentration of approximately $10^{17}/\text{cm}^3$ there. To effect good contact between the substrate and the FET bodies, it is also desirable to dope the graded-alloy region. Extensive work was done to achieve this. However, it was found that the use of diborane during the graded-layer formation reduced the deposition temperature at which surface pits, considered to be due to gas-phase nuclei, were formed. Additionally, chamber-wall deposition was seen to increase with the addition of diborane. Whether the reduction in film quality with the addition of diborane is an intrinsic effect or is due to nonidealities with the Stanford system is unclear. Nevertheless, it was decided to limit the boron doping to the surface region and use an "intrinsic" graded layer.

Chamber deposition was a considerable problem before it was recognized that the "standard" chamber cleaning procedure used with the reactor, which concluded with a 1050 C deposition of silicon on the susceptor with an atmospheric-pressure mixture of dichlorosilane and hydrogen, was "seeding" the chamber walls and facilitating the subsequent deposition of material there during the long epitaxial process. Replacement of this susceptor coat with a 750 C silane-and-hydrogen process substantially reduced the chamber-coating problem.

After device-grade epitaxy is achieved, the next challenge is the development of an insulator. The most promising candidate is probably an NO thermally grown germanium-oxinitride gate. However, NO, N_2O , and NH_4 atmospheric furnaces are not readily available in the Stanford lab, and therefore deposited SiO_2 was used, using silane and oxygen at 200 mtorr and 400 C. The problem with this method is that during the deposition thermal oxidation of the germanium at the surface can occur. An improvement in interface quality was observed when a thin silicon layer was deposited on top of the germanium immediately prior to oxide deposition. Immediately after etching away the field oxide over the device active area in dilute hydrofluoric acid, wafers are loaded into the reactor load-lock, purged in nitrogen for an hour, and loaded into the reaction chamber at an ambient temperature of 100 C or lower. After a further hydrogen purge, the temperature is ramped up to 400 C, where silane is flowed for 3 minutes, with hydrogen flow maintained throughout. Since the desorption of hydrogen from a silicon surface is generally the

rate-limiting step in silane CVD, and since hydrogen bonds much more readily with silicon than with germanium, this process is effectively self-limiting -- silicon deposits on the exposed germanium surface but, once the surface is all silicon, hydrogen bonds with the surface and the growth is virtually blocked. Oxide deposition immediately follows this process.

The gate electrode is also formed in the epitaxial reactor. *In-situ* boron-doped $\text{Ge}_{0.30}\text{Si}_{0.70}$ is readily deposited at 500 C with a resulting resistivity of 1 mohm-cm. No further activation anneal is required. Deposition is initiated with a silicon seed layer. This is made thick enough (at least several extrinsic Debye lengths) to establish a well-defined workfunction at the electrode-insulator interface. Then, to avoid problems associated with band discontinuities, the germanium fraction is gradually graded up to 30%. After the bulk of the gate is thus deposited, the germanium fraction is continuously reduced back to zero and the growth is completed with a silicon capping layer, used to present a well-understood surface for later processing.

The remaining fabrication is standard silicon MOS -- implant $10^{15}/\text{cm}^2$ arsenic at 25 keV, activate the dopant at 500 C, deposit an LTO sub-metal dielectric, etch contact holes, and deposit and pattern titanium and aluminum sputtered metal. Finally, a 275 C forming gas anneal is done to improve the oxide-semiconductor interface and the conductivity of the metal-semiconductor contacts.

Initial testing of completed devices is expected to begin by the end of March 1996. Testing of strained-Ge devices is expected in April.

Bibliography

- [Daembkes] Daembkes et al; *IEEE TED*, **33**:663 1986.
- [Fitzgerald] Fitzgerald et al; *APL* **59**:811-813 1991.
- [Garone] Garone et al; *IEEE EDL* **12**(5):230-232 1991.
- [Hymes] Hymes et al; *Journal of the Electrochemical Society*, **135**(4):961-965 1988.
- [Ismail] Ismail et al; *IEEE EDL* **13**(5):229-231 1992.
- [King] King et al; *IEEE EDL* **10**:52 1989.
- [LeGoues91] LeGoues et al; *Phys Rev Letters*, **66**(22):2903-2906 1991.
- [LeGoues92] LeGoues et al; *Journal of Applied Physics*, **71**:4230-4243 1992.
- [Meyerson] Meyerson et al; *Applied Phys Letters*, **53**(25):2555-2557 1988.
- [Nayak] Nayak et al; *IEEE EDL* **12**(4):154-156 1991.
- [Pearsall86] Pearsall and Bean; *IEEE EDL* **7**:308 1986.

- [Pearsall89] Pearsall; CRC Critical Reviews of in Solid State and Materials Sciences, **15**(6):551-600 1989.
- [Schaffler] Schaffler et al; *Semiconductor Device Tech*, **7**:260-266 1992.
- [Tsang] Tsang et al; *Appl Phys Let*, **62**(10):1146-1148 1993.
- [Welser] Welser et al; *IEDM Tech Digest* 1000-1002 1992.
- [Wang] Wang et al; *Materials Research Society Symp Proc*, **220**:403-408 1991.
- [Xie] Xie et al; *Materials Research Society Symp Proc*, **220**:413-417 1991.

UNIT: 6

TITLE: Portable Video on Demand in Wireless Communication

PRINCIPAL INVESTIGATOR: T. H. Meng

GRADUATE STUDENT: K. Precoda

I. Introduction

This research aims at providing low-power video compression for portable wireless video applications. We developed a power efficient video encoder architecture that uses pyramid vector quantization (PVQ) to compress video data. The decoded image quality using this encoder is better on average in terms of PSNR than JPEG.

In wireless communication, the available bandwidth generally changes with time. Our PVQ encoder, therefore, adjusts the frame rate according to the available bandwidth. If a large bandwidth is available, we increase the frame rate, improving the video quality at the receiver. If the bandwidth is limited, we decrease the frame rate, which results in degraded video quality. This ability to dynamically vary the compression rate allows the encoder to adaptively vary the amount of video data transmitted to achieve the best image quality for a given available bandwidth.

To handle variable frame rates while consuming the absolute minimal power, which is critical in portable systems, we propose to use circuits whose speed/power consumption can be adjusted by actual encoder throughput requirements. Our approach is to design a power supply controller that can adjust the DC voltage to control the desired performance. At high frame rates or when large bandwidth is available, the encoder would operate at high voltages, and, therefore, higher frequencies, allowing more image pixels to be processed per second. If smaller bandwidth is available, the supply voltage need not operate at a high voltage and is decreased appropriately to allow efficient operation at the required throughput. The encoder, therefore, consumes the absolute minimal power necessary to meet the frame rate of the encoder.

II. Power-Supply Regulation

In order to provide a variable supply voltage as a function of the processing speed required, the voltage regulator must rapidly vary the supply voltage to meet the required throughput rate, while maintaining high power efficiency. We have designed a dc-dc switching regulator that

achieves efficiency in excess of 90% with a tracking speed of under 1 ms. The regulator supplies efficiently from a few milli-Watts to several hundred milli-Watts for all supply voltages of interest.

A. Introduction to Switching Regulator

The switching regulator works by chopping the input battery voltage to generate a wave of pulses. These pulses pass through a second-order low-pass filter, which reduce the ac component to an acceptable ripple. The chopping is accomplished by active devices, which are integrated on a single chip to meet the size and weight requirements in portable applications. The inductor and capacitor, which form the low-pass filter, cannot be integrated to standard CMOS process, unfortunately, because of their large inductance and capacitance values. Consequently, off-chip inductors and capacitors are used.

B. Low Power Techniques For Switching Regulators

The switching regulator can ideally achieve 100% efficiency. There are three main sources of dissipation which cause the conversion efficiency to be less than unity: conduction loss in the chopping transistors, switching loss due to parasitics, and gate drive loss.

To improve the conversion efficiency, we employ synchronous rectification and fixed pulse-width voltage modulation. A diode is typically placed between a ground and the input to the low-pass filter to drive the pulse to zero volt. For low-power applications, the voltage drop across the diode causes significant power loss compared to the power delivered. This conduction loss is minimized by replacing a diode with a gated NMOS, which reduces the conduction loss substantially. This use of NMOS is called *synchronous rectification*.

The output voltage is approximately equal to the input voltage multiplied by the duty factor. The duty cycle can be changed arbitrarily by varying the pulse-width or keeping the pulse-width constant and varying the operation frequency. Unlike most traditional switching regulators, we use the latter approach of modulating the output voltage. By keeping the pulse-width constant and varying the operation frequency, the size of the optimal chopping transistors remains relatively constant for varying operating conditions. The amount of energy delivered per pulse remains invariant to varying load sizes, allowing a PMOS transistor sizing that is optimal for all loading conditions.

C. The Feedback Loop

The chopping of the supply voltage makes the converter intrinsically a nonlinear system. Methods of approximating this non linearity to a linear system for a small region of operation and

performing appropriate feedback compensation techniques are well known. Since our encoder must operate at wide load conditions as well as operating voltages, the location of the poles and zeros move by substantial amounts. To maintain stability with a fast response time, the converter needs to track the large movements of poles and zeros and place the compensating poles appropriately. This complicates the controller, which increases power dissipation and lowers efficiency. A nonlinear feedback controller is, therefore, employed requiring only a few adders and comparators. This controller is shown to be stable for all operating regions of interest.

III. Low-Power PVQ Encoder

A real-time, low-power video encoder for pyramid vector quantization is estimated to dissipate only 2.1mW for video compression of images of 256x256 pixels at 30 frames per second in 0.8-micron CMOS technology with a 1.5V supply. Applying this quantizer to subband decomposed images, the quantizer performs better than JPEG on average. We achieve this high level of power efficiency with image quality exceeding that of variable rate codes through efficient algorithm-to-hardware mapping.

Pyramid vector quantization (PVQ) is a quantization technique first introduced by Fisher as a fast method of quantizing and coding Laplacian-like data. PVQ is a fixed rate coding technique with compression performance asymptotically equivalent to a uniform scalar quantizer with entropy coding. PVQ compression capabilities and its fixed rate property, which prevents catastrophic error propagation, are well suited for transmission over noisy, wireless channels.

Our PVQ encoder operates by grouping multiple independent and identically distributed subband coefficients into a vector and finding the index of the nearest lattice point on the surface of a 1-dimensional hyperpyramid. The regularity of a pyramidal shape in multi-dimensional space allows simple recursive equations to assign unique indices to each point on the pyramid. Computing codeword indices rather than employing a large look-up table stored in memory is critical in low-power architectures, since memory operations consume far more power than arithmetic computations, often by several factors for on-chip memory accesses, not to mention the power needed for off-chip memory access.

The original PVQ encoder algorithm was modified to provide more parallelism and pipelining in the architecture, allowing the encoder to efficiently quantize vectors with dimensions as large as 256. This was the key to achieving high compression efficiency while maintaining good image quality at very low power levels.

IV. Conclusion

The goal of this research was to study the energy-on-demand design methodology for implementing low-power video compression systems. The methodology introduced using our dynamic variable supply voltage, however, can be employed in various other digital signal processing applications, where the required throughput rates are time-variant. We are exploring other applications for this energy-on-demand design methodology.

JSEP Supported Presentations

1. E. K. Tsern "A Low-Power Video-Rate Pyramid VQ Decoder," presentation at 1996 IEEE International Solid-State Circuits Conference, February 1996.
2. T. H.-Y. Meng, "A Low-Power Encoder Architecture for Pyramid Vector Quantization of 2-D Subband Coefficients," presentation at International Conference on Image Processing, October 1995.
3. N. Chaddha, "Scalable Video Compression," presentation at International Conference on Image Processing, October 1995.
4. T. H. Meng, "Portable Video-on-Demand in Wireless Communication," invited seminar presentation at Princeton University, April 1995.
5. T. H. Meng, "Portable Video-on-Demand in Wireless Communication," invited seminar presentation at Rockwell Science Center, March 1995.

JSEP Supported Publications

1. W. Namgoong, N. Chaddha and T. H. Meng, "Low-Power Video Encoder/Decoder Using Wavelet/TSVQ With Conditional Replenishment," *Proceedings IEEE ICASSP*, Atlanta, Georgia, May 1996.
2. E. K. Tsern and T. H. Meng, "A Low-Power Video-Rate Pyramid VQ Decoder," *1996 IEEE International Solid-State Circuits Conference Digest of Technical Papers*, pp. 162-163, San Francisco, California, February 1996.
3. W. Namgoong and T. H. Meng, "A Low-Power Video Encoder for Pyramid Vector Quantization of Subband coefficients," submitted to *Journal of VLSI Signal Processing*, February 1996.
4. A. C. Hung, E. K. T. and T. H. Meng, "Error Resilient Pyramid Vector Quantization for Image Compression," submitted to *IEEE Transactions on Image Compression*, January 1996.

5. B. M. Gordon, E. K. Tsern, and T. H.-Y. Meng, "Design of a Low-Power Video Decompression Chip Set for Portable Applications," invited submission to *Journal of VLSI Signal Processing*, October 1995.
6. W. Namgoong, M. Davenport, T. H.-Y. Meng, "A Low-Power Encoder Architecture for Pyramid Vector Quantization of 2-D Subband Coefficients," *Proceedings of 1995 IEEE Workshop on VLSI Signal Processing*, pp. 391-400, Osaka, Japan, October 1995.

UNIT: 7

TITLE: Adaptive DFE for GMSK in Indoor Radio Channels

PRINCIPAL INVESTIGATOR: J. M. Cioffi

GRADUATE STUDENTS: R. D. Wesel and K. Jacobsen

I. Introduction

Point-to-multipoint transmission problems are finding increasing application in broadcast and data communication networks. Such problems were the main focus of the supported JSEP research. Two Ph.D. students are matriculating in 1996 in these areas, Richard Wesel and Krista Jacobsen. Both have significant results, as reported below, and several published or pending papers under this contract's support.

Super-redundancy - R.D. Wesel

Rick Wesel's work focused on broadcast coding methods. In this area, a single source of digital information sends the same information to many remote users, with no feedback path. The transmission paths may vary from user to user and with time for a particular user. Such a situation is characteristic of terrestrial or satellite broadcast networks.

Rick found that to optimize a transmission system fully, the channel characteristic must be known to both the transmitter and the receiver. The consequent optimal action of the transmission system is then a function of this known channel characteristic. In the broadcast case, each user has a different channel characteristic and all are unknown to the transmitter. However, the maximum data rate that could be achieved by each of these users is roughly the same that should achieve at least the worst-case capacity on all the channels. Rick found this rate can be achieved without having to use different codes/designs for the different user paths.

Rick's work then progressed to a search for such a robust code, and several have been found as well as a general search procedure. These codes and the search procedure are described in Section II.

Multipoint-to-point access protocol and analysis - K. Jacobsen

The main focus of Krista Jacobsen's research has been the mechanisms for upstream access in a point-to-multipoint transmission architecture. The specific architecture studied was tree-structured coaxial networks, but the results also apply to wireless and local-area networks.

This work has produced a number of protocols and contention resolution methods for multicarrier transmission with such networks. In particular, a combination of time and frequency division access are combined at the physical transmission layer to improve throughput versus latency trade-offs in such networks, as described in Section III.

A method for network synchronization and coordination was postulated for a multicarrier transmission system and reservation-based access protocols were investigated. Significant improvements in throughput and efficiency were obtained with respect to time-only multiplexing.

Both sets of work have resulted in a reasonable level of publication as reported in Sections II. and III.

II. Trellis Codes for Correlated Fading - Rick Wesel

The Problem

Consider transmission over one or more channels subject to fading in time or frequency such that the fading can be estimated at the receiver but is unknown to the transmitter. An important example of this is digital video broadcast as shown in Fig. 1. This scenario also occurs in single carrier modulation in the presence of flat fading and frequency hopped transmissions in the presence of frequency selective fading. Our work during this period has resulted in a trellis code design technique that provides reliable performance on a much wider class of fading patterns than previous fading channel coded modulation techniques.

Performance of an example code designed using the new technique demonstrates consistently good bit error rate performance over a wide range of fading behavior. It is well known that an Ungerboeck--style trellis code concatenated with a forward error correction code can give performance approaching the channel capacity of an additive white Gaussian noise channel. The new code design technique combined with a forward error correction code can approach capacity simultaneously for a whole class of frequency selective channels, within the limitations of a fixed transmitter power spectrum. Thus, the new design technique provides trellis

codes that are ideal for use in broadcast transmissions where a single transmission must work for a variety of different channels.

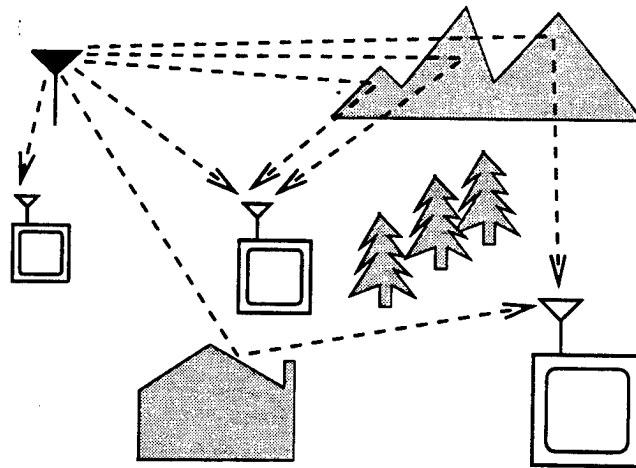


Figure 1: Digital Video Broadcast.

These new codes have already generated significant interest in industry. Telia Research, the Swedish telecommunications company, is exploring how these codes can be used to provide reliable wireless data links between a base station and a mobile user. Here again the transmitter cannot specialize the transmission to the particular fading. Unlike the broadcast situation, there is only one fading pattern. However, the transmitter does not know what that fading pattern is. Thus a robust code is required.

The Subchannel Decomposition

Regardless of whether the fading is in time or frequency, the overall subchannels with different SNRs are shown in Fig. 2. The new code design technique is based on two observations involving this subchannel decomposition. First, the number of coded bits transmitted per symbol needs to approach the subchannel capacities of the high capacity subchannels. Second, code distance must be carefully distributed to ensure that as many uncorrelated subchannels as possible can contribute to decoding.

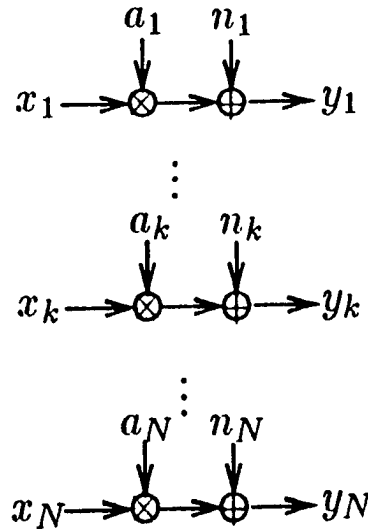


Figure 2: Overall subchannels with different SNRs.

Super-Redundancy

The first requirement, that the number of coded bits transmitted per symbol be large, implies that good fading--channel trellis codes will have a large amount of redundancy. This concept of super-redundancy can be contrasted with the additive white Gaussian noise channel, where it was shown that only one bit of redundancy is required [Ungerboeck]. In the fading environment, the subchannel capacities can vary by a large amount. To efficiently use the channel as a whole, each individual subchannel must be used efficiently. This requires that the number of coded bits be large enough that the high capacity subchannels can be fully utilized.

Code Distance Distribution and Correlation in Fading

It is well known that trellis codes designed for fading channels should distribute distance over as many codewords as possible. This concept of maximizing "effective code length" was formalized [Wilson and Leung] almost ten years ago. Our new code design technique takes this concept one step further. There is correlation between adjacent symbols in time or frequency, and our code design metrics consider the distance associated with groups of adjacent frequencies or time symbols together. This ensures that distance is not accidentally concentrated on a group of symbols which will experience correlated fading.

Previous techniques use interleaving to make the fading appear locally uncorrelated. This interleaving is combined with codes designed for fading assumed to be completely uncorrelated. Interleaving is still used in our new codes to allow short constraint lengths. However, unlike

previous techniques, the permuted correlation in the interleaved fading channel is a primary consideration the code design procedure.

To utilize this correlation information in a straightforward way, periodic interleaving is used. The interleaving period is chosen small enough that symbols within one period are essentially uncorrelated. Symbols separated by multiples of the interleaving period are extremely correlated. Thus symbol--error distances on symbols separated by multiples of the interleaving period provide exactly one "diversity branch".

The code design search procedure finds the trellis code that spreads code distance as evenly as possible on as many of these diversity branches as possible. The number of diversity branches in such a scheme is upper bounded by the period of the interleaver. However, if this period is chosen correctly, that is also the limit of the diversity present in the fading environment. Detailed discussions of the code design procedure can be found in the publications listed at the end of this section.

Performance of the New Codes

To see how well the new codes can perform we consider the example of multicarrier broadcast and consider the four different frequency responses shown in Fig. 3. A multicarrier system with 512 subcarriers is assumed and the desired information rate will be fixed at 1 bit per symbol. Our code design procedure produces a rate 1/4 convolutional code which is used to select points from a 16 QAM constellation. This code is compared with a standard code for multicarrier broadcast of 1 bit per symbol -- a rate 1/2 code used to select points from a 4 PSK constellation. Both codes have 64 states and thus require Viterbi decoders with the same complexity.

Figure 4 shows that the newly designed code provides consistent performance on all four of these channels. At a bit error rate of 10^6 the new code has all four performance curves within a band of 0.75 dB. The standard code performs 1 dB better on the Flat Channel (Channel 1). However, its performance becomes unacceptable as the frequency selectivity becomes more pronounced. On the Step Channel (Channel 4), which is a step in the frequency response, the standard code has bit error rates close to 1/2 for the entire range of the plot.

Conclusion

The new codes produced by this research provide reliable performance over a wide variety of time/frequency fading patterns. This type of consistent reliability is unmatched by previous

techniques, and the new codes will find applications in numerous data communication applications including digital video broadcasting and wireless data networks.

III. Design and Analysis of Multipoint-to-point Discrete-Multitone-based Networks - Krista S. Jacobsen

The Problem

As the deployment of hybrid fiber-coax (HFC) networks by both cable television and telephone companies continues, efficient, cost-effective techniques to transmit digital multimedia signals both to and from the home must be developed. Transmission channels in the downstream direction (from the central site to the customer premise) are generally high-quality, and use of a single-carrier modulation in broadcast mode is probably sufficient for downstream transmission. However, the upstream bandwidth of HFC networks is often plagued by numerous transmission impairments, including passband ripple, spectral nulls, and radio-frequency ingress. Hence, a robust upstream modulation technique is required to ensure that effective communications can occur in the presence of these impairments. Furthermore, because HFC networks are generally configured in tree-and-branch topologies, as shown in Fig. 1, the return channel (upstream bandwidth) is shared among many users, potentially thousands. Consequently, use of the available upstream bandwidth must be coordinated somehow to ensure the channel is used efficiently.

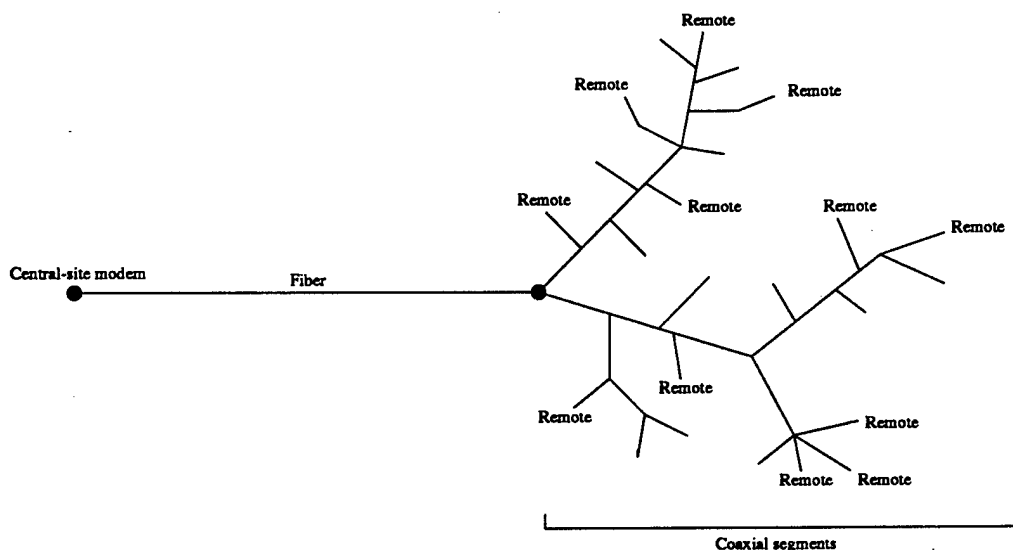


Figure 1: HFC network configuration.

Discrete-multitone (DMT), a type of multicarrier modulation, has been shown in previous JSEP-sponsored papers (i.e., [Jacobsen and Cioffi-a], [Jacobsen and Cioffi-b]) to offer significant advantages for upstream transmission in HFC channels, particularly because DMT can optimize the data rate on channels suffering from severe transmission impairments such as those prevalent in the HFC return channel. However, to exploit these advantages and achieve improvements in the overall network performance, an efficient multicarrier-based channel access protocol is required. Many general protocol alternatives can be adapted to control multicarrier-based remote units, but the challenge lies in developing a simple, easily managed protocol that separates, as much as possible, the modulation from the protocol without losing the benefits gained by using multicarrier transmission.

The goal of this work is to design and analyze DMT-based multipoint-to-point networks in general and HFC networks in particular. Issues addressed include the gains in spectral efficiency that can be realized by using DMT modulation rather than a single-carrier modulation with equalization; methods for installing, synchronizing, and training multiple DMT-based remote units [Jacobsen and Cioffi-c]; the performances of various channel access protocol alternatives for DMT-based networks; and the differences in overall system complexity of single-carrier-based and DMT-based HFC networks. In this document, recent work on the installation, synchronization, and training procedures is presented along with preliminary results on one channel access protocol designed specifically for multicarrier-based multipoint-to-point networks.

Procedure for Installing, Synchronizing and Training DMT-based Remote Units

Because HFC networks span large distances and must support a large number of remote units, at all times, including during installation, it is desired that communications from any particular remote unit to the central-site modem occur without disrupting in-progress data transmissions by other remote units. Thus, when it is first installed on the network, before a remote unit is synchronized it may only transmit while the other remote units are silent. First, the remote unit loop-times its local clock with the central-site model clock, which is broadcast in a downstream control channel. To enable remote unit installations, synchronization, and initial training, silent intervals of a predetermined length are observed periodically in the upstream data stream by all remote units on the network. The central unit transmitter constantly sends a trigger in a downstream control channel to instruct installing remote units to send their installation parameters. Upon receiving a valid installation signal during the silent period, the central-site controller compares the signal's symbol boundaries to those of symbols transmitted by the remote units currently operating. In general, there will be a difference in the symbol boundaries, and the controller computes and sends the time delay required of the synchronizing unit to correct the

misalignment. The remote unit then implements the requested sample delay and transmits a signal requesting verification that it is synchronized. If the remote unit transmission is indeed synchronous, the central unit controller sends a signal to that unit in the downstream channel to indicate that no further shifting is required, and that the remote unit may now communicate with the central-site modem incorporating the appropriate delay. Otherwise, the synchronization procedure is repeated until the central-site controller determines the remote unit is synchronized. After the initial symbol delay has been determined, unless a remote unit is moved or its connection to the network is terminated, it should not have to be resynchronized. Failing to synchronize the remote units to within a certain tolerance can result in interchannel interference, which can decrease the achievable bit rates on the affected subchannels.

After receiving and incorporating the required sample delay from the central-site modem, an installing remote unit transmits a wide-band signal during a specified number of upcoming silent periods to train the central unit receiver. Because the newly installed remote unit is now synchronized with respect to the other remote units, it can transmit using all of the symbols during the next several silent periods for channel analysis. All other remote units remain quiet while the remote unit transmits a training signal on the permissible subset of the subchannels allocated to it, and the central unit controller records the bit capacity and magnitude and phase of each subchannel from that remote unit. The bit capacities are used to determine subchannel assignments when the remote later requests either a constant data rate or a packet transmission. Because the controller allocates the subchannels to the various remote units every symbol period, it can apply the appropriate subchannel magnitude/phase inverse to each subchannel to demodulate the received signal. Hence, if the remotes are all properly synchronized, the signal arriving at the central unit receiver, which is actually an aggregate of transmissions from a number of different remote units, can be demodulated as though it were from a single remote modem, using the appropriate mixture of subchannel magnitude/phase inverses.

After a remote has been installed, it is periodically retrained during another silent interval reserved specifically for this purpose. As during the installation silent period, all remote units that are not training remain quiet to allow the central unit controller to update its settings for the training remote. Depending on the frequency of these silent intervals, the number of remotes on a particular network, and other system parameters, each remote could be retrained as often as many times per second or as infrequently as every few seconds.

Design and Analysis of the Reservation-Based Multicarrier (RBM) Protocol

After the remote units have been installed, synchronized, and trained, they are capable of transmitting without interfering with other remote units as long as they obey a channel access protocol. One alternative for controlling transmissions from remote units so that data is always transmitted collision-free is a reservation-based protocol. Under a generalized reservation-based protocol, to obtain permission to transmit data a remote unit must first transmit a reservation request. When a reservation has been granted, then the corresponding data message is guaranteed to be received intact (channel noise notwithstanding) by the central-site modem. If reservation requests are transmitted using the same bandwidth as data transmissions, then coordination of reservation requests is necessary to ensure they do not interfere with data transmissions.

The *Reservation-Based Multicarrier (RBM)* [Jacobsen and Cioffi-d] protocol has been developed for multicarrier-based multipoint-to-point networks (such as HFC) in which data transmissions scheduled by a central controller are desirable because remote units are unable to detect whether or not the upstream channel is in use, and data transmissions and reservation requests occupy the same bandwidth. Under the RBM protocol, each multicarrier symbol is marked by the central controller as either reserved for data transmissions or available for reservation requests. The controller, which resides at the central-site, broadcasts in a downstream control channel a binary-valued "channel status" signal during every symbol period to inform the remote units whether the channel will be in use for data transmissions during the subsequent symbol period. Based on the status of the next symbol, a remote unit with a data message ready for transmission either reschedules the request for a later time according to a nonpersistent algorithm, or it sends its reservation request during the next symbol period on a randomly-selected "frequency-domain slot." A frequency-domain slot is a set of subchannels in the frequency domain that together support, during one symbol period, a bit capacity equal to the number of bits required to transmit a "request for bandwidth" (RFB). Each RFB consists of at least the remote unit's address, and it may also indicate the size of forthcoming data messages, quality of service parameters, etc, depending on the network particulars. As an example of frequency-domain slot partitioning, if there are N B -bit subchannels in the multicarrier system, and RFBs consist of kb bits, then the subchannels are grouped into K sets of k subchannels. Because all remote units with data ready for transmission may send their RFBs during any symbol period not reserved for scheduled data transmissions, RFBs are subject to collisions with other RFBs. However, because RFBs are generally short, the collision probability is fairly low, and partitioning the subchannels into K sets ensures that collisions always overlap completely. A system allowing different numbers of bits on each subchannel may group a different number of subchannels into each of the K sets, but collisions still overlap completely, and the concept of frequency-domain slots still

holds. Regardless of what method is used to divide the subchannels into frequency-slots, the partitioning must be observed by all remote units.

After transmitting its RFB on one of the K subchannel sets, a remote unit waits a specified period of time, determined by the round-trip propagation delay of the channel and the central unit processing time, to ascertain whether or not its RFB arrived successfully at the receiver. If the waiting remote does not receive a grant message from the central controller within a certain period of time, which indicates that its RFB collided with another unit's RFB or was unintelligible to the receiver for some other reason, it reschedules the RFB for a later time according to a delay distribution. If the remote does receive a grant message before timing out, it begins to transmit its message using all subchannels during the symbol period corresponding to the index sent by the central controller. Figure 2 illustrates the protocol timing, channel status signal, and upstream channel activity when a successful RFB occurs and the minimum delay is incurred.

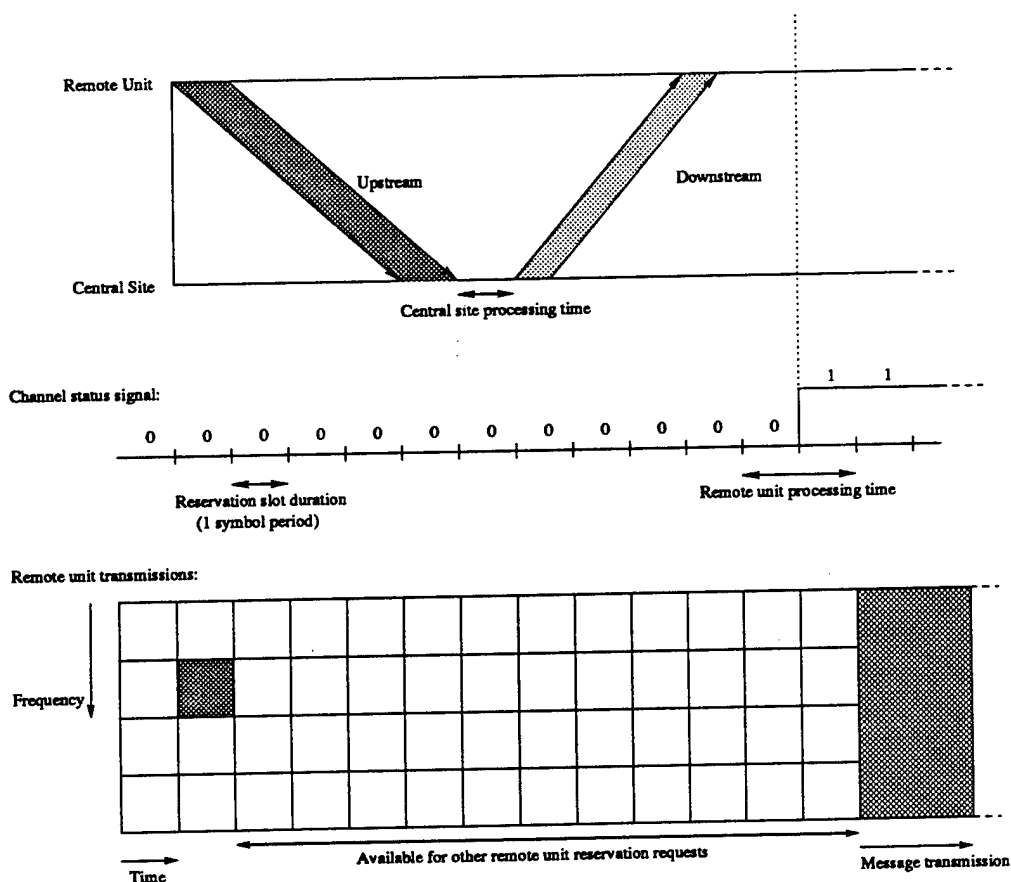


Figure 2: Illustration of protocol with $K = 4$.

To simplify protocol management, all multicarrier remote units are constrained to transmit using the same bit tables. In other words, for all remote units, the number of bits b_i on the i th subchannel is the same. Note that the number of bits supported by subchannel i need not equal the number supported by subchannel j , as long as b_i and b_j are the same across all remote units on the network. Under the constraint of equivalent bit tables, the central unit receiver applies the same decoding procedure to every received symbol. Therefore, the receiver does not need to know in advance which of the remotes is transmitting an RFB or, for that matter, a message. Furthermore,

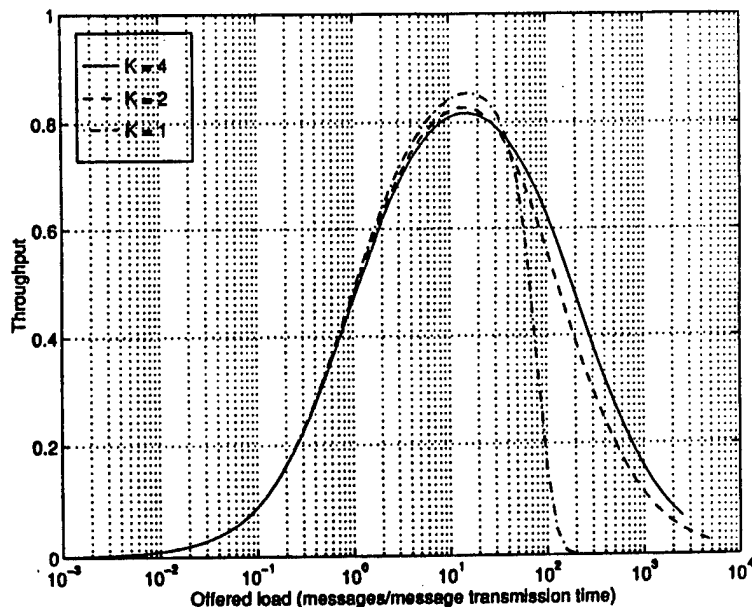


Figure 3: Achievable throughputs of 32-subchannel ($K = 4$), 16 subchannel ($K = 2$), and 8 subchannel ($K = 1$, slotted single carrier) RBM protocols: 256 bits/message, 2 bits/subchannel, 16 bits/RFB.

although RFBs are transmitted on only a few of the available subchannels, messages are transmitted using all subchannels in a symbol. Hence, if a symbol has been reserved for data transmission, only one remote unit transmits during that symbol. This restriction simplifies protocol management and enables the remote units to determine whether to transmit RFBs by checking only the binary channel status signal.

Of interest for evaluating the performance of a particular protocol is the expected *throughput* it enables at various loads. The throughput is defined as the percent of time the channel is used for data transmissions. Figure 3 shows the achievable throughputs of the RBM protocol with various numbers of subchannels as a function of offered load for a network with $a = 1$. A constant message length of 256 bits was used with 16 bits per RFB, and the available transmission

bandwidth was divided into 32 ($K = 4$), 16 ($K = 2$), and 8 ($K = 1$, slotted single-carrier) 2-bit subchannels. Hence, the time required to transmit each message is the same for each scenario, and the achievable throughputs for the various values of K may be compared without modification.

The figure illustrates that the throughput achieved by the RBM protocol is a function of the offered traffic load and the expected number of successful RFB transmissions during a message slot, which in turn is a function of the number of subchannel sets, K , available for simultaneous reservation requests. For loads in the range from $G = 0$ to $G = G_C$, where G_C is some critical value of the offered load, the slotted single-carrier version of the RBM protocol ($K = 1$) can achieve a slightly higher throughput than $K > 1$ versions, whereas the $K > 1$ versions perform significantly better when $G > G_C$. As a general rule, then, the throughput of networks with higher numbers of subchannels degrades less severely as the offered load increases. It is important to note that the offered load in a real system is likely to change significantly during operation. Hence, at times the network load will exceed G_C , especially if the remote unit traffic sources are bursty. The results presented imply that certain networks intended to support bursty traffic sources can benefit significantly from using multicarrier modulation in conjunction with the RBM protocol.

References

- [Ungerboeck] G. Ungerboeck, "Channel Coding with Multilevel/Phase Signals," *IEEE Transactions on Information Theory*, 40(5):1459-1473, September 1994.
- [Wilson and Leung] S. G. Wilson and Y. S. Leung, "Trellis-Coded Phase Modulation on Rayleigh Channels," *Proceedings of ICC-87*, pages 739-742, June 1987.
- [Jacobsen and Cioffi-a] K. S. Jacobsen and J. M. Cioffi, "An Efficient Digital Modulation Scheme for Multimedia Transmission on the Cable Television Network," in *Technical Papers, 43rd Annual National Cable Television Association (NCTA) Convention and Exposition*, New Orleans, LA, May 1994.
- [Jacobsen and Cioffi-b] K. S. Jacobsen and J. M. Cioffi, "High-performance Multimedia Transmissions on the Cable Television Network," in *Proceedings 1994 International Conference on Communications*, New Orleans, LA, May 1994.
- [Jacobsen and Cioffi-c] K. S. Jacobsen and J. M. Cioffi, "Synchronized DMT for Multipoint-to-point Communications on HFC Networks," in *Globecom '95 Proceedings*, Singapore, November 1995.

[Jacobsen and Cioffi-d] K. S. Jacobsen and J. M. Cioffi, "Achievable Throughput in Multicarrier-based Multipoint-to-point Networks Using a Reservation-based Channel Access Protocol," submitted to *Globecom '96*.

JSEP Supported Publications

1. R. D. Wesel and J. M. Cioffi, "Fundamentals of Coding for Broadcast OFDM," In *Proceedings of the 29th Asilomar Conference on Signals Systems & Computers*, November 1995.
2. R. D. Wesel and J. M. Cioffi, "A Transmission System Using Codes Designed for Transmission with Periodic Interleaving," U.S. Patent Pending.
3. R. D. Wesel and J. M. Cioffi, "Trellis Codes for Channels with Correlated Fading," in Preparation for Submission to *IEEE Transactions on Communications*.
4. K. S. Jacobsen and J. M. Cioffi, "An Efficient Digital Modulation Scheme for Multimedia Transmission on the Cable Television Network," in *Technical Papers, 43rd Annual National Cable Television Association (NCTA) Convention and Exposition*, New Orleans, LA , May 1994.
5. K. S. Jacobsen and J. M. Cioffi, "High-performance Multimedia Transmissions on the Cable Television Network," in *Proceedings 1994 International Conference on Communications*, New Orleans, LA, May 1994.
6. K. S. Jacobsen and J. M. Cioffi, "Synchronized DMT for Multipoint-to-point Communications on HFC Networks," in *Globecom '95 Proc.*, Singapore, Nov. 1995.
7. K. S. Jacobsen and J. M. Cioffi, "Achievable Throughput in Multicarrier-based Multipoint-to-point Networks Using a Reservation-based Channel Access Protocol," submitted to *Globecom '96*.

TITLE: Robust Estimation Methods for Adaptive Filtering**PRINCIPAL INVESTIGATOR: T. Kailath****GRADUATE STUDENTS: Y. C. Pati and B. Hassibi****1 Introduction**

Our earlier JSEP-supported work was concerned with the use of spatial and temporal (signal) structure in smart antennas for mobile radio networks. The work done there gradually led us to consider, and to study, the *robustness* of the underlying algorithms with respect to model uncertainties and lack of statistical information. In particular, of interest were adaptive filtering algorithms which are widely used in communications (as well as in many other areas) for the identification and equalization of channels.

Classical methods for such problems require a priori knowledge of the statistical properties of the signals. In many applications, however, one is faced with model uncertainties and lack of statistical information. Therefore the aforementioned methods are not directly applicable. Moreover, it is not even clear what the behaviour of such estimation schemes might be if the assumptions on the statistics and distributions are not exactly met.

Adaptive filtering techniques are currently widely used to cope with such model uncertainties and lack of a priori knowledge. The methods currently used fall into the two general classes of least-squares-based algorithms (such as recursive-least-squares or RLS) and gradient-based algorithms (such as least-mean-squares or LMS). While the former class is derived from an explicit cost function, it is suspect whether their robustness properties are always desirable. On the other hand, the former methods are rather ad-hoc and do not follow from a rigorous framework. However, the gradient algorithms are by far the ones most used in applications. Our work now provides some analytic explanation of this fact.

In the last decade such problems have received great attention in control theory, where a so-called H^∞ approach has been extensively studied. It turns out, in particular, that the LMS algorithm is H^∞ -optimal, thus establishing the observed robustness of this very widely used algorithm. We have also obtained some results on the robustness of least-squares-based adaptive filters. This framework is currently being used to explore new adaptive filtering algorithms for nonstationary scenarios.

2 Adaptive Filtering

The standard model assumed in adaptive filtering is the following:

$$d_i = h_i^T w + v_i, \quad i \geq 0 \quad (1)$$

where $\{d_i\}$ is an observed output sequence (often called the reference signal), $\{h_i\}$ is a known input vector sequence, w is an unknown weight vector that we intend to estimate, and $\{v_i\}$ is an unknown disturbance, which may also include modeling errors. We shall make no assumptions on the statistics or distributions of the $\{v_i\}$.

We denote the estimate of the weight vector using all the information available up to time i by

$$w_i = \mathcal{K}(d_0, d_1, \dots, d_i; h_0, h_1, \dots, h_i).$$

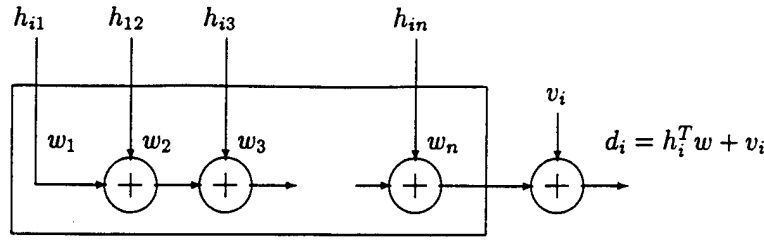


Figure 1: The model for adaptive filtering.

2.1 Least-Squares-Based Methods

There are a variety of choices for w_i , but the most widely used estimate w_i , is one that satisfies the following least-squares (or H^2 criterion):

$$\min_w \left[\mu^{-1} |w - w_{-1}|^2 + \sum_{j=0}^i |d_j - h_j^T w|^2 \right], \quad (2)$$

where w_{-1} is the initial estimate of w , and $\mu > 0$ represents the relative weight that we give to our initial estimate compared to the “sum of squared-error” term $\sum_{j=0}^i |d_j - h_j^T w|^2$. In the so-called *pure least-squares* problems one takes $\mu = \infty$, so that the first term in the cost function of (2) does not appear.

The *exact* solution to the above criterion is the RLS (Recursive Least Squares) algorithm:

$$w_i = w_{i-1} + k_{p,i}(d_i - h_i^T w_{i-1}), \quad w_{-1} \quad (3)$$

with $k_{p,i} = \frac{P_i h_i}{1 + h_i^T P_i h_i}$ and $P_{i+1} = P_i - \frac{P_i h_i h_i^T P_i}{1 + h_i^T P_i h_i}$, $P_0 = \mu I$.

RLS has certain stochastic optimality properties: if we assume in model (1) that the $w - w_{-1}$ and $\{v_i\}$ are zero mean independent Gaussian random variables with variances μI and 1 respectively, then the RLS algorithm yields the maximum likelihood estimate of w_i . In particular, it minimizes the expected *prediction error energy*:

$$E \|e_p\|_2^2 = E \sum_{j=0}^i |h_j^T w - h_j^T w_{j-1}|^2. \quad (4)$$

2.2 Gradient-Based Algorithms

In gradient-based algorithms instead of exactly solving the least-squares problem (2), the estimates of the weight vector are updated along the negative direction of the *instantaneous* gradient of the cost function appearing in (2). Two examples are the LMS (Least-Mean-Squares)

$$w_i = w_{i-1} + \mu h_i (d_i - h_i^T w_{i-1}), \quad w_{-1} \quad (5)$$

and normalized LMS

$$w_i = w_{i-1} + \frac{\mu}{1 + \mu h_i^T h_i} h_i (d_i - h_i^T w_{i-1}), \quad w_{-1} \quad (6)$$

algorithms. Note that in the case of LMS the gain vector $k_{p,i}$ in RLS (which had to be computed by propagating a Riccati equation) has been simply replaced by μh_i . Likewise if we compare normalized LMS with the RLS algorithm, we see that the difference is that instead of propagating the matrix P_i via the Riccati recursion we have simply set $P_i = \mu I$, for all i . Therefore the LMS and normalized LMS algorithms were long considered to be *approximate* least-squares solutions and were thought to lack a rigorous basis.

2.3 The Question of Robustness

We noted that under suitable stochastic assumptions, H^2 -optimal adaptive filters have certain desirable optimality properties. However, a question that begs itself is what the performance of such filters will be if the assumptions on the disturbances are violated, or if there are modelling errors in our model so that the disturbances must include the modelling errors? In other words

- *is it possible that small disturbances and modelling errors may lead to large estimation errors?*

Obviously, a nonrobust algorithm would be one for which the above is true, and a robust algorithm would be one for which small disturbances lead to small estimation errors.

The problem of robust estimation is thus an important one. As we shall see in the next section, the H^∞ robust estimation formulation is an *attempt* at addressing this question. The idea is to come up with estimators that minimize (or in the suboptimal case, bound) the maximum energy gain from the disturbances to the estimation errors. This will guarantee that if the disturbances are small (in energy) then the estimation errors will be as small as possible (in energy), *no matter what the disturbances are*. In other words the maximum energy gain is minimized over *all possible* disturbances. The robustness of the H^∞ estimators arises from this fact. Since they make no assumption about the disturbances, they have to accomodate for all conceivable disturbances, and are thus over-conservative. So this is not necessarily the best solution and we are also exploring weaker criteria.

3 The H^∞ Approach

We now apply the H^∞ methodology to adaptive filtering (see [Hassibia] for details).

To this end, consider the following three types of estimation errors:

(i) The prediction error:

$$e_{p,i} = h_i^T w - h_i^T w_{i-1}.$$

(ii) The filtered error:

$$e_{f,i} = h_i^T w - h_i^T w_i.$$

(iii) The smoothed error:

$$e_{s,i} = h_i^T w - h_i^T w_N.$$

Any choice of estimation strategy \mathcal{K} will induce a transfer operator from the disturbances to the above estimation errors. These we shall denote by

(i) $T_p(\mathcal{K})$: transfer operator from the disturbances $\{w - w_{-1}, v_i\}$ to the prediction errors $\{e_{p,i}\}$.

(ii) $T_f(\mathcal{K})$: transfer operator from the disturbances $\{w - w_{-1}, v_i\}$ to the filtered errors $\{e_{f,i}\}$.

(iii) $T_s(\mathcal{K})$: transfer operator from the disturbances $\{w - w_{-1}, v_i\}$ to the smoothed errors $\{e_{s,i}\}$.

Now for any choice of estimator, \mathcal{K} , and any given input disturbance sequence $\{w - w_{-1}, v_i\}$, that yields the prediction error, $\{e_{p,i}\}$, we may compute the energy gain

$$\frac{\|e_p\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2}, \quad (7)$$

where $\|a\|_2^2 = \sum_j |a_j|^2$ is defined as the energy of the sequence $\{a\}$ and μ is a positive constant. Thus (7) is a measure of the "amplification" of the noise given our choice of estimator \mathcal{K} . Now the ratio in (7) will clearly depend on the input disturbance, $\{w - w_{-1}, v_i\}$. however, if we consider *all* possible disturbance sequences $\{w - w_{-1}, v_i\}$, then we can find the largest energy gain in (7). This leads us to the definition of the H^∞ norm of a transfer operator T .

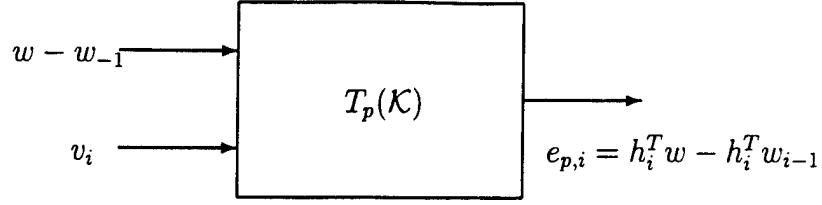


Figure 2: Transfer operator from the unknown disturbances $\{w - w_{-1}, v_i\}$ to the prediction errors $\{e_{p,i}\}$. Likewise for $T_f(\mathcal{K})$ and $T_s(\mathcal{K})$.

Definition 1 (The H^∞ Norm) The H^∞ norm of a transfer operator T is defined as

$$\|T\|_\infty = \sup_{x \in h^2, x \neq 0} \frac{\|Tx\|_2}{\|x\|_2} \quad (8)$$

where h^2 denotes the space of all square-summable causal sequences.

We now propose to choose the estimator \mathcal{K} so as to minimize the H^∞ norms of $T_p(\mathcal{K})$, $T_f(\mathcal{K})$ and $T_s(\mathcal{K})$. To be more specific we have the following problem.

Problem 1 (H^∞ Adaptive Filtering Problem) Find estimators $w_i = \mathcal{K}_p(d_0, \dots, d_i; h_0, \dots, h_i)$, that minimize the maximum energy gain from disturbances to estimation errors for each of the aforementioned errors, i.e., find estimation strategies \mathcal{K}_p , \mathcal{K}_f and \mathcal{K}_s such that

$$\gamma_p^2 = \inf_{\mathcal{K}_p} \sup_{w, v \in h_2} \frac{\|e_p\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2} \quad (9)$$

$$\gamma_f^2 = \inf_{\mathcal{K}_f} \sup_{w, v \in h_2} \frac{\|e_f\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2} \quad (10)$$

and

$$\gamma_s^2 = \inf_{\mathcal{K}_s} \sup_{w, v \in h_2} \frac{\|e_s\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2} \quad (11)$$

where $|w - w_{-1}|^2 = (w - w_{-1})^T(w - w_{-1})$ and $\mu > 0$ reflects a priori knowledge of how close w_{-1} is to w .

It turns out that nice solutions can be obtained for all three problems. The solutions to Prob. 1 are given below (see [Hassibia]), in which we have assumed that the input vectors $\{h_i\}$ are such that

$$\lim_{N \rightarrow \infty} \sum_{i=0}^N h_i^T h_i = \infty.$$

Solution to (i): If μ satisfies the bound

$$0 < \mu < \inf_i \frac{1}{h_i^T h_i} \quad (12)$$

then $\|T_p(\mathcal{K})\|_\infty$ is minimized by the LMS algorithm with learning rate μ ,

$$w_i = w_{i-1} + \mu h_i (d_i - h_i^T w_{i-1}), \quad w_{-1}$$

and the minimum H^∞ norm is given by

$$\gamma_p = 1.$$

Remarks:

- (a) The fact that $\gamma_p = 1$ indicates that there is no amplification of the disturbances. Thus the prediction error energy will never exceed the disturbance energy.
- (b) The above result is true only if the learning rate μ satisfies the bound (12). This is in accordance with the well-known fact that LMS behaves poorly if the learning rate is chosen too large.

Solution to (ii): $\|T_f(K)\|_\infty$ is minimized by the normalized LMS algorithm

$$w_i = w_{i-1} + \frac{\mu}{1 + \mu h_i^T h_i} h_i (d_i - h_i^T w_{i-1}) \quad , \quad w_{-1}$$

and the minimum H^∞ norm is given by

$$\gamma_f = 1.$$

Remark: Note once more that there is no amplification of the noise. Now, however, we have no restriction on μ .

Solution to (iii): $\|T_s(K)\|_\infty$ is minimized by the least-squares solution, and the minimum H^∞ norm is

$$\gamma_s = 1.$$

Remark: Thus least-squares algorithms are H^∞ optimal with respect to smoothing errors.

4 Robustness of Least-Squares Algorithms

Now that we have developed the H^∞ optimality of the LMS and normalized LMS algorithms with respect to prediction and filtered errors, it is natural to ask what the performance of the RLS algorithm will be with respect to these error criteria.

In order to answer the above question we need to compute the H^∞ norm of the RLS algorithm. Finding this H^∞ norm essentially amounts to finding the maximum singular value of a linear time-varying operator. Upper bounds on the H^∞ norm can be found by checking for the positivity of the solution of a certain time-varying discrete-time Riccati recursion. Although both approaches can be used in principle, they require knowledge of *all* the input data vectors $\{h_i\}$.

Since in adaptive filtering problems we are given, and are forced to process, the data in real time, we cannot store all the data and use the aforementioned methods to compute bounds for the H^∞ norm. Therefore the main effort in the results given below is to obtain bounds on H^∞ norm that use simple a priori knowledge of the $\{h_i\}$ and not their explicit values [Hassibib].

(i) For RLS, we can show

$$(\sqrt{R} - 1)^2 \leq \sup_{w, v \in h_2} \frac{\|e_p\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2} \leq (\sqrt{R} + 1)^2$$

or to give a "looser" bound

$$(\sqrt{1 + \mu \bar{h}^2} - 1)^2 \leq \sup_{w, v \in h_2} \frac{\|e_p\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2} \leq (\sqrt{1 + \mu \bar{h}^2} + 1)^2,$$

where

$$R = \max_i 1 + h_i^T P_i h_i, \quad \bar{h}^2 = \max_i |h_i|^2, \quad \underline{h}^2 = \min_i |h_i|^2.$$

Remark: Note that for large μ , the H^∞ norm grows as μ . This shows that the pure least-squares problem (with $\mu = \infty$) is highly non-robust with respect to prediction errors.

(ii) For filtered errors we have

$$\sup_{w, v \in h_2} \frac{\|e_f\|_2^2}{\mu^{-1}|w - w_{-1}|^2 + \|v\|_2^2} \leq (\sqrt{1/r} + 1)^2 \leq 4,$$

where

$$r = \min_i 1 + h_i^T P_i h_i \geq 1.$$

Remarks:

- (a) Note that, as with normalized LMS, the H^∞ norm does not depend on μ .
- (b) The above result for filtered errors is an intermediate stage between the smoothing error case (where the H^∞ and H^2 optimal filters coincide) and the prediction error case (where the performance of LMS and RLS can be drastically different.)

5 Future Work

The H^∞ approach to adaptive filtering described in the previous section suggests several directions for future research. We mention a few here.

5.1 Time-Varying Problems

So far we have assumed that the weight vector, w , is constant in time. In many applications one needs to assume a time-varying, w , and must therefore devise algorithms that can track the time-variations of the weight vector.

In such cases, one approach is to use windowing. Two common windowing schemes are the following.

- (i) **Exponential Window:** The exponential window gives (exponentially) larger weight to the more recent data. In particular, the prediction error and disturbance energies are computed as:

$$\sum_{j=0}^i \lambda^{-j} |e_j|^2 \quad \text{and} \quad \sum_{j=0}^i \lambda^{-j} |v_j|^2, \quad (13)$$

where $0 < \lambda < 1$ is the so-called forgetting factor that is chosen based upon a priori knowledge of how fast the weight vector varies with time.

- (ii) **Finite-Memory Window:** In this case one only considers the last L data points so that the prediction error and disturbance energies are computed as

$$\sum_{j=i-L+1}^i |e_j|^2 \quad \text{and} \quad \sum_{j=i-L+1}^i |v_j|^2, \quad \text{respectively.} \quad (14)$$

L is often referred to as the window length.

It is therefore useful to consider the H^∞ filters that result from such “windowed” definitions of energy. The filters that are obtained in this fashion will have good tracking properties and, at the same time, be robust.

5.2 Mixed H^2/H^∞ Estimation

Fig. 5.2 shows the (squared) singular values of $\mathcal{T}_{p,rls}$ and $\mathcal{T}_{p,lms}$ (the transfer operators from disturbances to estimation errors for RLS and LMS) for $N = 50$ (where N is the number of observed data points) and $\mu = .9$, for a simple one-dimensional adaptive filtering problem. As can be seen the maximum singular value for $\mathcal{T}_{p,lms}$ is one, whereas for $\mathcal{T}_{p,rls}$ it is much larger. On the other hand, the RLS algorithm minimizes the Frobenius norm (the sum of the squared singular values) of the transfer operator \mathcal{T}_K which can be visualized as the area under the curve of the (squared) singular values. Thus if we choose disturbances uniformly from the space C^{50} , the RLS algorithm will have better average performance than LMS, although its worst-case performance is significantly worst.

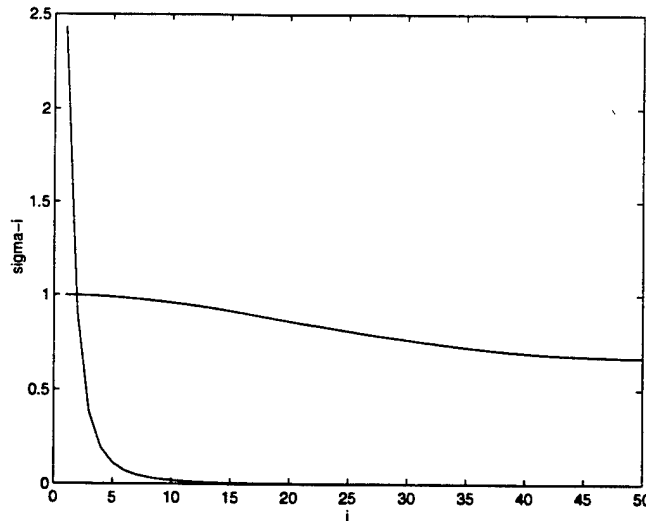


Figure 3: Singular values for $\mathcal{T}_{p,rls}$ and $\mathcal{T}_{p,lms}$ for $N = 50$ and $\mu = .9$.

Note, moreover, that although the LMS algorithm does not allow any amplification of the disturbances, it does not provide significant suppression of the disturbances, either. (The smallest squared singular value for $\mathcal{T}_{p,lms}$ which represents the minimum energy gain is roughly 0.65.) Since the H^∞ optimal filters are not unique (LMS is only the central solution), it is very interesting to study the possibility of choosing other H^∞ optimal filters to further reduce the Frobenius norm of \mathcal{T}_K . This will result in algorithms that have the best possible average behaviour while at the same time having the best possible worst-case performance. This framework is called the mixed H^2/H^∞ estimation framework and is an area that we intend to pursue.

References

- [Hassibia] B. Hassibi, A.H. Sayed and T. Kailath. H^∞ Optimality of the LMS Algorithm. *IEEE Trans. on Signal Processing*, vol 44, pp. 267-281, February 1996.
- [Hassibib] B. Hassibi and T. Kailath. H^∞ bounds for the recursive-least-squares algorithm. in *Proceedings of the 33rd IEEE Conference on Decision and Control*, pp. 3927-3929, Orlando, FL, Dec 1994.

TITLE: Efficient Data Compression**PRINCIPAL INVESTIGATOR: T. Cover****GRADUATE STUDENTS: E. Erkip, P. Fahn, G. Iyengar,
D. Kimber, I. Kontoyiannis, A. Lapidoth and E. Ordentlich****AFFILIATED STUDENTS: V. Castelli and T. Jing**

1 Scientific Objectives

We apply techniques of information theory to problems of efficiently compressing data, for the purpose both of storage and transmission of data. Our results concern compression schemes both for specific circumstances, such as specific noise profiles in communications channels, and for general circumstances, such as human "non-algorithmic" image compression. How much can be gained by customizing data compression to specific circumstances, as opposed to using the convenience of established, off-the-shelf algorithms? During the past year this work resulted in 9 supported papers and 1 Ph.D. thesis.

2 Summary of Research

One of the new innovations in this year's research on data compression is the work of G. Iyengar on the capacity of the voice channel. In most channels, one constructs waveforms that are invulnerable to noise. In the voice channel, we use an independent source of white noise to drive a filter which is chosen at the discretion of the speaker. The filter models the speaker's vocal tract. We ask how many distinguishable filters there are in the presence of additive white noise as heard by the listener. A solution of this problem will yield the capacity of the voice channel, a characterization of an optimal vocabulary, the role of feedback in speaking, models for languages, and finally, since we're able to model speech, the optimal method of data compression associated with the voice channel.

The results of a multiyear project on image compression are being summarized in a paper by T. Jing, T. Cover and R. Wesel and L. An. We have found that if humans are allowed to intervene in a legitimate image compression experiment substantial improved data compression can be achieved over existing JPEG standards. This is not a surprise because our experiment requires over 100 man hours of work per image to achieve the desired data compression. However, the magnitude of the improvement shows how far existing image compression algorithms have to go before they achieve their ultimate limits.

In other work by P. Fahn, we are investigating the quantum theoretic correlations in distributed measurements and observations. This leads naturally to questions of data compression for quantum theoretic systems, or by means of such systems. These data compression results will turn out to be nonclassical because the marginal distributions of measurements have been shown, via Bell's inequality, to be inconsistent with any multivariate distribution. In a sense, there is no underlying physical reality with which the observations are consistent.

3 Detailed Research Descriptions

3.1 Image Compression

Our experiment to compare the image compression abilities of humans and computers is in its final stages. Our goal is to estimate the minimal rate, in bits per pixel, at which an image can be compressed without incurring significant perceptible distortion. First, one experimental subject simplifies a given image without significantly distorting it, and then another subject predicts the simplified image, pixel by pixel, as accurately as possible. The accuracy of the second subject's predictions can be quantified to yield an estimate of the entropy of the simplified image. Not only will our results be useful as a benchmark to researchers in the field, but the experimental framework itself may lead to a new algorithm for data compression. A paper detailing the results of the experiment is currently under preparation.

3.2 Voice Channel

The thrust of this research is to develop a characterization of the capacity and optimal coding vocabularies of voice channels, which are mathematical models intended to capture properties of human speech generation. This research area will provide guidance on data compression for a voice channel or other channels with similar characteristics.

Consider a communication system with a channel characterized by a linear filter g in an additive Gaussian noise environment, i.e., $y(t) = u(t) * g(t) + z(t)$, where $z(t)$ is white noise and $u(t)$ is the channel input. Instead of fixing the filter $g(t)$, which is the traditional approach, we fix the input signal $u(t)$ and attempt to choose a distribution on the space of linear, passive, causal filters $g(t)$ that maximizes the mutual information between the output and the filter. This model and its discrete-time analog are, we propose, an approximate model for the voice generation process

3.3 Feedback in Communication

It was recently shown by [Pombra and Cover] that the maximum achievable throughput (sum of rates of all users) of a Gaussian multiple access channel with feedback is at most twice that achievable without feedback. We prove [Ordentlich] a somewhat stronger result which establishes the factor of two bound not only for the total throughput but for the entire capacity region as well. Specifically, we show that the capacity region of a Gaussian multiple access channel with feedback is contained within twice the capacity region without feedback.

We have recently extended the factor of two bound on the capacity region of Gaussian multiple access channels to channels with inter-symbol interference (ISI). For single user Gaussian channels there is no information theoretic complication introduced by the addition of a causal linear filter at the transmitter. If the filter is invertible, the channel can be transformed into an ISI-free channel with an appropriately modified noise spectrum. For the multiple access channel, if the ISI filters are not identical for all transmitters, as is the case in practice, no such transformation is possible. This new result demonstrates that in wireless communications networks, once steady state has been reached via power control and channel learning, the maximum additional gain in capacity region afforded by receiver-to-transmitter feedback is limited to a factor of two, no matter how cleverly the feedback is used.

3.4 Robustness of Communication

Lapidoth, in a series of papers [Lapidoth 1], [Lapidoth 2], [Lapidoth 3], has considered the robustness of signaling in the presence of noise in an unknown environment. It is well known that Gaussian signals and matched filter decoding is optimal for signaling with a power constraint over an additive white noise Gaussian channel. This is the basis for much of the signaling which is done, say, in deep space communication or in mobile communication. Lapidoth is able to show that even if you fix the receiver to be a matched filter receiver and continue to use the same signals, the information will get through the channel no matter what the noise is, just so long as the total noise power is not increased. Specifically, if the distribution of the noise is changed from Gaussian and independent to non-Gaussian and arbitrarily time dependent, as long as the noise power is not increased, the channel will still work and the probability of error will be exponentially small. This shows the robustness of existing communication schemes to changes in the underlying assumptions on the model.

3.5 Universal lossless compression

We investigate the convergence properties of optimal data compression algorithms. The two major directions of this research are entropy (redundancy) estimation for text and images, and calculation of the exact distribution of codeword-lengths for large data sets.

In a series of papers [Kontoyiannis and Suhov 1], [Kontoyiannis and Suhov 2], [Kontoyiannis, Algoet and Suhov], we have investigated the convergence properties of several entropy estimation algorithms. These algorithms are suggested by optimal data compression schemes that are based on pattern-matching, such as the celebrated Lempel-Ziv algorithm for text compression. Our results prove the optimality of some existing methods and also suggest new algorithms for the efficient estimation of the redundancy withing any given data set.

The results along the second direction of this research provide a second-order analysis of the distribution of the size of losslessly encoded data. We prove a second-order refinement to Shannon's (lossless) Source Coding Theorem. In essence this result says that the distribution of the deviations of the compressed data size from its mean is, at best, Gaussian. The minimum variance of this Gaussian is a quantity characteristic of the source, and it provides a theoretical bound on the variance of the encoded data size.

3.6 Quantum Data Compression

As feature sizes on electronic devices continue to shrink, there is fear that quantum effects will become significant and hinder reliable functioning. A few researchers have started to realize, however, that the strange properties of quantum systems can be exploited to great advantage. A quantum bit is not restricted to the "classical" values 0 and 1, but can take values anywhere in a two-dimensional Hilbert space. If we can find practical methods to compress data using quantum bits, therefore, the efficiency achieved could greatly surpass the limits of classical data compression.

Of course if quantum bits are so efficient, they may be used throughout the computation and communication processes. Therefore we have been studying the possibilities for quantum data compression of quantum data as well as of classical data. If a quantum source generates data described by a density operator ρ , data compression limits can be established using the Von Neumann entropy $-k\text{Tr}\rho \log \rho$, instead of the Shannon entropy used in traditional data

compression. Furthermore, one can compress partly entangled pairs of quantum particles into a small number of completely entangled pairs — the so-called Bell states — which can then be used for efficient communication of quantum data.

References

- [Castelli and Cover] V. Castelli and T. Cover. On the Exponential Value of Labeled Samples. *Pattern Recognition Letters*, 16:105-111, January 1995.
- [Cover and King] T. Cover and R. King. A Convergent Gambling Estimate of the Entropy of English. *IEEE Trans. on Information Theory*, IT-24(4):413-421, July 1978.
- [Kontoyiannis and Suhov 1] I. Kontoyiannis and Yu. M. Suhov, Prefixes and the Entropy Rate for Long-range Sources, in *Probability Statistics and Optimization*, F. P. Kelly, ed. Chichester, England: Wiley, 1994, pp. 89-98.
- [Kontoyiannis and Suhov 2] I. Kontoyiannis and Yu. M. Suhov, Stationary Entropy Estimation via String Matching, in *Proceedings of the Data Compression Conference*, Snowbird, Utah, April 1996.
- [Kontoyiannis, Algoet and Suhov] I. Kontoyiannis, P.H. Algoet and Yu. M. Suhov, Two Consistent Entropy Estimates for Stationary Processes and Random Fields, in preparation to be submitted to *IEEE Transactions on Information Theory*.
- [Lapidoth 1] S. Shamai and Amos Lapidoth. Bounds on the Capacity of a Spectrally Constrained Poisson Channel. *IEEE Transactions on Information Theory*, 39(1):19-29, January 1993.
- [Lapidoth 2] Amos Lapidoth. On the Reliability Function of the Ideal Poisson Channel with Noiseless Feedback. *IEEE Transactions on Information Theory*, 39(2):491-503, March 1993.
- [Lapidoth 3] Amos Lapidoth. The Performance of Convolutional Codes on the Block Erasure Channel Using Various Finite Interleaving Techniques. *IEEE Transactions on Information Theory*, 40(5):1459-1473, September 1994.
- [Ordentlich] E. Ordentlich. A Class of Optimal Coding Schemes for Moving Average Additive Gaussian Noise Channels with Feedback. *Proceedings of the IEEE International Symposium on Information Theory*, p.467, June 1994.
- [Pombra and Cover] S. Pombra and T. Cover. Non-White Gaussian Multiple Access Channels with Feedback, *IEEE Transactions on Information Theory*, 40(3):885-892, May 1994.

4 Publications Supported by JSEP

4.1 Ph.D. Theses Supported by JSEP

A. Lapidoth, "Mismatched Decoding of the Multiple-Access Channel and Some Related Issues in Lossy Source Compression," August 1995.

4.2 Published Papers Supported by JSEP

1. S. Pombra and T. Cover. Non-White Gaussian Multiple Access Channels with Feedback, *IEEE Transactions on Information Theory*, 40(3):885-892, May 1994.
2. Z. Zhang and T. Cover. On the Maximum Entropy of the Sum of Two Dependent Random Variables. *IEEE Transactions on Information Theory*, 40(4):1244-1246, July 1994.
3. V. Castelli and T. Cover. On the Exponential Value of Labeled Samples. *Pattern Recognition Letters*, 16:105-111, January 1995.

4.3 Papers Submitted for Publication

1. T. Cover and E. Ordentlich. Universal Portfolios with Side Information. To appear in *IEEE Transactions on Information Theory*.
2. P. Fahn. Maxwell's Demon and the Entropy Cost of Information. To appear in *Foundations of Physics*.
3. I. Kontoyiannis, P.H. Algoet and Yu. M. Suhov, Two Consistent Entropy Estimates for Stationary Processes and Random Fields, in preparation to be submitted to *IEEE Transactions on Information Theory*.
4. A. Lapidoth. Mismatched Decoding and the Multiple-Access Channel. *Stanford University Statistics Department Technical Report No.87*, February 1995. Under review by *IEEE Transactions on Information Theory*.
5. E. Ordentlich. On the Factor-of-Two Bound for Gaussian Multiple Access Channels with Feedback. Submitted to *IEEE Transactions on Information Theory*.
6. E. Ordentlich and T.M. Cover. Max-Min Optimal Investing. To appear in *Proceedings of 1996 IEEE/IAFE Conference on Computational Intelligence for Financial Engineering*.